

University of Southampton Research Repository ePrints Soton

Copyright © and Moral Rights for this thesis are retained by the author and/or other copyright owners. A copy can be downloaded for personal non-commercial research or study, without prior permission or charge. This thesis cannot be reproduced or quoted extensively from without first obtaining permission in writing from the copyright holder/s. The content must not be changed in any way or sold commercially in any format or medium without the formal permission of the copyright holders.

When referring to this work, full bibliographic details including the author, title, awarding institution and date of the thesis must be given e.g.

AUTHOR (year of submission) "Full thesis title", University of Southampton, name of the University School or Department, PhD Thesis, pagination

UNIVERSITY OF SOUTHAMPTON

**Agent-Based Trust and Reputation in
the Context of Inaccurate Information
Sources**

by

W. T. Luke Teacy

A thesis submitted in partial fulfillment for the
degree of Doctor of Philosophy

in the

Faculty of Engineering, Science and Mathematics
School of Electronics and Computer Science

December 2006

UNIVERSITY OF SOUTHAMPTON

ABSTRACT

FACULTY OF ENGINEERING, SCIENCE AND MATHEMATICS
SCHOOL OF ELECTRONICS AND COMPUTER SCIENCE

Doctor of Philosophy

by **W. T. Luke Teacy**

Trust is a prevalent concept in human society that, in essence, concerns our reliance on the actions of other entities within our environment. For example, we may rely on our car starting to get to work on time, and on our fellow drivers, so that we may get there safely. For similar reasons, trust is becoming increasingly important in computing, as systems, such as the Grid, require integration of computing resources, across organisational boundaries. In this context, the reliability of resources in one organisation cannot be assumed from the point of view of another, as certain resources may fail more often than others. For this reason, we argue that software systems must be able to assess the reliability of different resources, so that they may choose which of them to rely on.

With this in mind, our goal is to develop mechanisms, or models, to aid decision making by an autonomous agent (the truster), when the consequences of its decisions depend on the actions of other agents (the trustees). To achieve this, we have developed a probabilistic framework for assessing trust based on a trustee's past behaviour, which we have instantiated through the creation of two novel trust models (TRAVOS and TRAVOS-C). These facilitate decision making in two different contexts with regard to trustee behaviour. First, using TRAVOS, a truster can make decisions in contexts where a trustee can only act in one of two ways: either it can cooperate, acting to the truster's advantage; or it can defect, thereby acting against the truster's interests. Second, using TRAVOS-C, a truster can make decisions about trustees that can act in a continuous range of ways, for example, taking into account the delivery time of a service.

These models share an ability to account for observations of a trustee's behaviour, made either directly by the truster, or by a third party (reputation source). In the latter case, both models can cope with third party information that is unreliable, either because the sender is lying, or because it has a different world view. In addition, TRAVOS-C can assess a trustee for which there is little or no direct or reported experience, using information about other agents that share characteristics with the trustee. This is achieved using a probabilistic mechanism, which automatically accounts for the amount of correlation observed between agents' behaviour, in a truster's environment.

Contents

Nomenclature	xiii
Acknowledgements	xvii
1 Introduction	1
1.1 The Meaning of Trust	2
1.2 The Relevance of Trust in Computing	4
1.3 Service-Oriented Computing	5
1.4 Research Objectives	7
1.5 Research Contributions	10
1.6 Thesis Structure	12
2 Computational Models of Trust	15
2.1 The Cognitive Viewpoint	15
2.2 Learning from Direct Observations	17
2.2.1 Improvised Models of Trust	17
2.2.2 Probabilistic Models of Trust	21
2.2.3 Dempster-Shafer Models of Trust	23
2.3 Learning from Others	25
2.3.1 External Factors	26
2.3.2 The Majority Opinion	27
2.3.3 Past Performance	29
2.4 Mechanism Design	30
2.5 Assessing Unknown Agents	31
2.6 Summary	32
3 A Probabilistic Framework for Modelling Trust & Reputation	37
3.1 Background	38
3.2 Basic Notation and Problem Definition	42
3.3 Trust Assessment Based on Direct Observations	44
3.4 Reputation Framework	45
3.5 A Word of Warning on Sufficiency	49
3.6 Coping with Inaccurate Reputation	50
3.6.1 Statistical Noise	50
3.6.2 Opinion View	51
3.6.3 Opinion Incentives	51
3.7 Summary	53

4	TRAVOS: A Trust Model for Boolean Action Spaces	55
4.1	Instantiating the Framework for Boolean Action Spaces	56
4.2	Instantiating the Reputation Mechanism	59
4.3	Filtering Inaccurate Reputation	61
4.3.1	Estimating the Probability of Accuracy	62
4.3.2	Adjusting Reputation Source Opinions	63
4.4	Reputation Gathering for TRAVOS	65
4.4.1	Reputation Brokering	65
4.4.2	When to Seek Reputation	67
4.5	An Application to Agent-Based Virtual Organisations	68
4.5.1	System Overview	68
4.5.2	Walk-through Scenario	71
4.5.2.1	Calculating Trust	72
4.5.2.2	Calculating Reputation	73
4.5.2.3	Handling Inaccurate Opinions	74
4.6	Empirical Study	75
4.6.1	Experiment Methodology	76
4.6.2	TRAVOS Against the Beta Reputation System	77
4.6.3	TRAVOS Component Performance	79
4.7	Summary	79
5	TRAVOS-C: A Trust Model for Continuous Action Spaces	83
5.1	The TRAVOS-C Model	84
5.2	Learning Group Behaviour	86
5.3	Instantiating TRAVOS-C for Continuous Action Spaces	89
5.3.1	The Parameter Domains	89
5.3.2	The Reputation Function	89
5.3.3	Parameter Distributions for Group Behaviour	92
5.4	Applying the Model	93
5.5	A Monte-Carlo Method for TRAVOS-C	96
5.5.1	Hyperparameter Sampling	97
5.5.1.1	Conditional distribution for β	99
5.5.1.2	Conditional distribution for α	99
5.5.1.3	Conditional distribution of m	100
5.5.1.4	Conditional distribution of v	100
5.5.2	Parameter Sampling	101
5.5.2.1	Conditional distribution of μ	103
5.5.2.2	Conditional distribution of τ	104
5.5.3	Sampling Methods for Conditional Parameter Distributions	105
5.6	Empirical Study	109
5.6.1	Experiment Methodology	109
5.6.2	Basic Learning Behaviour	112
5.6.3	Learning from Reputation	115
5.6.4	Learning Reputation Source Correlations	117
5.6.5	Performance under Assumption Violations	119
5.7	Summary	123

6	Conclusions and Future Work	125
6.1	Thesis Summary	125
6.2	Research Contributions	127
6.2.1	Communicating Reputation	128
6.2.2	Addressing Inaccurate Reputation	128
6.2.3	Assessing Trust based on Group Behaviour	129
6.3	Limitations	130
6.4	Future Work	131
6.4.1	Dynamic Behaviour	131
6.4.2	Correlation Between Tasks	131
6.4.3	Implications of Reputation in Group Learning	132
6.5	Conclusions	132
A	Techniques for Numerical Integration	135
A.1	Deterministic Methods	135
A.2	Monte Carlo Methods	137
A.3	Independent Samples and Rejection Sampling	139
A.4	The Metropolis-Hastings Algorithm	142
A.5	Gibbs Sampling	144
A.6	Slice Sampling	144
B	Parameter Mapping for the Beta Distribution	147
	Bibliography	151

List of Figures

2.1	The three way relationship between trust, fear and authority.	16
3.1	Venn diagram of overlapping reputation datasets.	45
4.1	Example beta pdf plots; note that when $\alpha = 1, \beta = 1$ (top-left) the distribution is uniform in the interval $[0, 1]$	58
4.2	Example beta distributions for aggregating opinions of 3 agents.	60
4.3	Illustration of $\rho_{a_{tr}, a_{rep}}$ estimation process.	63
4.4	Reputation brokering system.	66
4.5	The CONOISE-G architecture.	69
4.6	TRAVOS reputation system versus BRS.	78
4.7	TRAVOS component performance.	80
5.1	The TRAVOS-C model.	85
5.2	Bayesian network for inferring group priors.	86
5.3	Examples of group parameter distributions with behaviour samples.	87
5.4	Examples of gamcon type II, precision and gamma distributions for comparison.	106
5.5	Example gamcon type II densities with Gaussian and Cauchy approximates.	107
5.6	Cauchy proposal densities for precision distributions.	108
5.7	Parameter estimates with variance sum of 25, varying direct observations.	113
5.8	Parameter estimates with variance sum of 9, varying reputation observations.	114
5.9	Behaviour parameter estimates, based on reliable reputation.	116
5.10	Behaviour parameter estimates, based on unreliable reputation.	117
5.11	Reputation noise parameter estimates with evidence for correlation.	118
5.12	Example skewed distributions, generated using transformed gamma densities.	120
5.13	Example bimodal distributions, generated using a mixture of two Gaussian densities.	121
5.14	Behaviour parameter estimates, based on reports from a lying reputation source.	122
A.1	An example of an integrand, approximated with different numbers of rectangles.	136
A.2	Example Laplacian function approximations.	137
A.3	Example rejection sampling regime.	140
A.4	Examples of poor rejection regimes.	141
A.5	Illustration of the slice sampling algorithm.	145

List of Tables

2.1	Trust value semantics used by Abdul-Rahman <i>et al.</i>	18
2.2	Frequencies of successful and unsuccessful interactions with different agents.	22
4.1	Combination of beta distributions.	64
4.2	Agent a_1 's interaction history with phone call service provider agents.	72
4.3	Agent a_1 's interaction history with HTML content service provider agents.	72
4.4	Agent a_1 's calculated trust and associated confidence level for HTML content and phone call service provider agents.	73
4.5	Agent a_1 's adjusted values for opinions provided by a_8 , a_9 and a_{10}	75
4.6	Observations made by a_1 given opinions from reputation sources. m represents that the interaction (to which the opinion applied) was successful, and likewise n means unsuccessful.	75
4.7	Reputation source populations.	77
5.1	Parameter set definitions.	96

List of Algorithms

1	Reputation broker update algorithm, performed by reputation sources. . .	67
2	The TRAVOS-C Gibbs sampler.	97
3	The TRAVOS-C simulation algorithm.	110
4	The rejection sampling algorithm.	140
5	The Metropolis-Hastings algorithm.	142
6	The Gibbs sampling algorithm.	143

Nomenclature

\propto	the proportional-to operator
$A - B$	The difference between sets A and B , defined as the set $\{x x \in A \wedge x \notin B\}$
\therefore	the therefore symbol
\emptyset	the empty set
\mathbb{R}	the set of real numbers
\mathbb{R}^+	the set of positive real numbers
\mathbb{Z}	the set of integers
\mathbb{N}	the set of natural numbers
$E[f(x)]$	expected value of a function $f(x)$
$p(x)$	probability density function (p.d.f.) or probability function (p.f.) of x
μ	mean parameter
σ^2	variance parameter
\mathcal{A}	the set of agents
a_i	the i th agent
a_{tr}	the truster
a_{te}	the trustee
a_{rep}	a reputation source
\mathcal{G}	the set of all distinct groups of agents
G_i	the i th group of agents
$G_{a_{tr}}$	the group of agents containing a_{tr}
t'	the current time step
$\mathcal{O}^{\mathcal{C}}$	the set of possible interaction outcomes in a context \mathcal{C}
$O_{a_{tr}, a_{te}}$	the outcome of an interaction with a_{te} , from the point of view of a_{tr}
$O_{a_{tr}, a_{te}}^t$	the outcome an interaction with a_{te} at time t , from the point of view of a_{tr}
$O_{a_{tr}, a_{te}}^{t:t+n}$	the set of outcomes between a_{te} and a_{tr} that occurred between time t and $t+n$
$\tilde{O}_{a_{rep}, a_{te}}$	an outcome $O_{a_{rep}, a_{te}}$, with added noise
$\tilde{O}_{a_{rep}, a_{te}}^{t:t+n}$	the set of outcomes $O_{a_{rep}, a_{te}}^{t:t+n}$, with noise added to each of its members
$\Theta^{\mathcal{C}}$	the set of distribution parameter vectors in context \mathcal{C}
θ	a distribution parameter vector
$\theta_{a_{tr}, a_{te}}$	the parameter vector for a_{te} 's behaviour distribution toward a_{tr}
$\epsilon_{a_{rep}}$	the parameter vector for the noise distribution associated a_{rep}
$\mathcal{N}^{\mathcal{C}}$	the set of possible noise values in context \mathcal{C}

$N_{a_{rep}}$	a noise term added to an outcome by a_{rep}
$\Phi^{\mathcal{C}}$	the set of hyperparameters for context \mathcal{C}
ϕ	a hyperparameter vector
ϕ_{G_i}	the hyperparameter vector associated with group G_i
$\phi_{a_{tr}, a_{te}}$	the hyperparameter vector for the distribution of $\theta_{a_{tr}, a_{te}}$
$R_{a_{rep}, a_{te}}$	the opinion of a_{rep} about a_{te}
$\mathcal{H}_{a_{tr}, a_{rep}}$	the history of opinions given to a_{tr} by a_{rep}
$\rho_{a_{tr}, a_{rep}}$	the estimated accuracy of a_{rep} according to a_{tr}
$\gamma_{a_{tr}, a_{te}}$	the confidence of a_{tr} in its assessment of a_{te}

Acknowledgements

The work presented in this document is part of the CONOISE-G project, funded by the DTI and EPSRC through the Welsh e-Science Centre, in collaboration with the Office of the Chief Technologist of BT. Many thanks go to my supervisors, Nick Jennings and Michael Luck, for their guidance and keen eye for detail; my colleagues on the CONOISE-G team, especially Jigar Patel my partner in trust, and Nir Oren for always knowing when something is missing. Finally, I wish to thank my parents, who above all else have made me the person I am today, for their sense of justice, for their love, and teaching me never to give up.

If I have seen further it is by standing on the shoulders of giants.

— Issac Newton (1642–1727), Letter to Robert Hooke, February 5, 1675

*To Michelle, my love, for putting her dreams on hold so that I
could fulfill mine.*

Chapter 1

Introduction

In human society, the fulfilment of even our most basic needs and desires is constantly at the mercy of other people's actions. That is, whether we are concerned with having food on our table or our post delivered on time, someone other than ourselves will play a critical role in making those events possible. Unfortunately, there is often a great deal of uncertainty surrounding the behaviour of our fellow beings: as we cannot in general read minds, so we cannot be certain about other people's intentions; likewise, we cannot always tell if the people we rely on have sufficient competence and resources to fulfill their obligations.

Managing this uncertainty is something we do almost subconsciously in our daily lives. For instance, when we need to delegate a task in the workplace, we normally choose a person who we believe is willing and able to do the job in hand (unless perhaps we have no better option). Also, we may choose not to disclose information to someone if we believe they will use that information to our disadvantage. Both of these cases, and many more, involve assessing the future action of a person or other entity, and deciding how we personally should act in response to that assessment. In such cases, it is common to talk about the notion of *trust*.

The concept of trust is thus prevalent in society and we use it in many contexts. As with many words in natural language, it is a term that is used frequently, understood implicitly, but not well defined. For now, we shall defer discussion on a more precise definition. However, we observe from this that trust can be associated with scenarios where action is required on behalf of some person or other entity, when the fulfilment of that action is not necessarily a foregone conclusion. This is particularly true when cooperation between different people is involved.

Increasingly, however, the need to deal with such scenarios is now also widespread in computer science. For example, in e-commerce, money regularly changes hands between

individuals and organisations that have no physical involvement; and in computer security, users must be trusted not to abuse their access rights, and those who are not trusted to perform certain tasks must be prevented from doing so.

Recent trends in computing, in particular the move towards large-scale open systems, threatens to make these types of concerns all the more challenging. Visions such as the semantic web and grid computing aim to enable the integrated use of computer resources across both geographical and organisational boundaries. It is envisaged that these systems will have highly dynamic properties, in which decisions regarding what resources to use, or allow access to, must be made rapidly, and perhaps automatically, in response to changing circumstances.

Managing the issues of trust that arise in these systems is a challenging problem and one that, in part, we aim to address in this thesis. To do so, we shall present a number of mechanisms that can be used to support automated decision making in the context of trust, but before we can do so, we need to be clear about what it is we mean by trust, and the precise set of problems that we aim to address.

To this end, the rest of this chapter sets the scene for our work by outlining our aims and objectives, along with a more precise definition of the types of problems we wish to address. Specifically, we break this discussion into six sections: Section 1.1 discusses the meaning of trust, and in particular what we mean by the term in this thesis; Section 1.2 outlines some of the general areas in computing for which trust is relevant; Section 1.3 looks at issues of trust in the particular area of service oriented computing, which this research is targeted at; Sections 1.4 and 1.5 detail the specific objectives and contributions of the thesis; and finally, Section 1.6 describes the structure of the rest of the thesis.

1.1 The Meaning of Trust

In our discussion so far, we have identified the notion of trust with scenarios that involve decision making in which the actions of different entities are relevant, but by no means certain. However, trust as an explicit concept is not one that has a single accepted definition. For example, definitions given by the Oxford English dictionary include, “Confidence in or reliance on some quality or attribute of a person or thing, or the truth of a statement.” which emphasises trust as both a belief and a dependence on someone or something. From a social science perspective, [Misztal \(1996\)](#) gives an in depth discussion of many different aspects of trust, and surveys a number of definitions that emphasize its role as both degree of belief, social relationship or acceptance without proof or investigation.

For our purposes, this is not a debate that we need enter into. Nevertheless, we are interested in a set of problems for which the term trust can meaningfully be applied, and it is useful to give a precise definition of the term, as it is used throughout the rest of this thesis. In doing so, we adopt the following definition, adapted from [Gambetta \(1988\)](#), which we believe captures the notion of trust we are interested in.

“Trust is a particular level of the subjective probability with which an agent will perform a particular action, both before she can monitor such an action, and in a context in which it affects her own action.”

There are five points in this definition that warrant further elaboration, and which we describe below:

1. **Trust between a pair of entities** — Trust is an assessment of one entity, which we shall refer to as the *trustee*, from the perspective of another entity, which we shall refer to as the *truster*. Although, we may occasionally refer to trust in oneself, normally these entities are distinct.
2. **Trust relates to a particular action** — Although sometimes we talk generally about our trust in an individual, a high level of trust in someone to perform one type of action does not imply a high level of trust in them to perform another. For example, just because we can trust a person to pick up a pen does not mean we trust them to run the country!
3. **Trust is a subjective probability** — Trust is subjective, because it is assessed from the unique perspective of the truster. It is dependent both on the individual set of evidence available to the truster and her relationship with the trustee.
4. **Trust is defined to exist before the respective action can be monitored** — Trust is a prior belief about an entity’s actions. It is an assessment made in a context of uncertainty. Once the truster knows the outcome of an action, she no longer needs to assess trust in relation to that outcome. Consider the difference in the statements, ‘I know you have brushed your teeth’ and ‘I trust that you have brushed your teeth’.
5. **Trust is situated in a context in which it affects the truster’s own action** — By this, we mean that our interest is limited only to those actions of a trustee that have relevance to the truster. Specifically, we are interested in trustee actions that, if their outcomes are known, would usefully inform the truster’s action decisions.

In our context, this is a strong definition because it captures both the purpose of trust that we are interested in, and its nature in a form that can be reasoned about analytically.

That purpose is to aid an entity to make decisions, with the entity in our case possibly being an automated computer system. As such, having an explicit concept that can be reasoned about is a prerequisite, if it is to provide a meaningful label to anything that can be automated by a computer. Defining trust as a subjective probability fulfills this requirement, because this is already a term that is well defined and understood within mathematics.

1.2 The Relevance of Trust in Computing

Having a clear definition of trust is one step to understanding its relevance to computing. However, there is a board range of issues in computing for which the term trust has previously been applied. Not all of these issues can be addressed appropriately in the same way, and it is beyond the scope of our work to do so. Nevertheless, to understand our work in context, it is necessary to give an overview of the types of issues encountered.

In general, issues of trust arise in computing when users and software can interact with information services, computing resources and other users with whom they are unfamiliar or have no physical contact. For instance, we may ask if we trust an information service to provide us with accurate information, or a particular website to respect the privacy of credit card details. The participation of large numbers of entities with conflicting interests in a large open system means that these examples are not isolated. In particular, we identify three important (possibly overlapping) areas in computing that are concerned with trust:

Security — Broadly, computer system security can be viewed as an attempt to limit the actions that individuals or software can perform with a given computer system. We can view this problem as reverse trust or, equivalently, fear (see Section 2.1): trust is generally concerned with a wish for an entity to perform an action, whereas fear is concerned with a wish for an entity *not* to perform an action. In this respect, we wish to avoid malicious actions, such as manipulation of important data or reading of trade secrets, and therefore an attempt to limit the ability to perform such actions only to those who are unlikely to have incentive to act maliciously.

Traditionally, computer security has been concerned with lower level issues such as: authentication, whereby the identity of a user is determined; authorisation, in which access to resources is granted; and data encryption (Gollmann, 1998; Pfleeger, 2002). Recently however, some in this field have started to refer explicitly to issues of trust, though in some cases, the term has been used merely as a synonym for authorisation or authentication (Grandison and Sloman, 2000). Others refer to it as a richer concept, and see it as a prerequisite condition for

authorisation. In this vein, Blaze et al. (1996) introduce the concept of *trust management*, which is concerned with specifying and applying *security policies* that state precisely what actions can be performed by a given entity.

Service Provision — In contrast to security, service provision concerns actions that a trustee is obliged to perform. Prime examples of this can be found in the semantic web (Berners-Lee et al., 2006), pervasive computing (Adelstein et al., 2004) and the Grid (Foster and Kesselman, 2004), in which certain tasks may be automatically delegated to systems that are outside the truster’s direct control. In this context, there may be a number of competing systems that can fulfill a particular task, each providing a different quality of service. Obviously, it is in the best interest of the truster to delegate in a way that maximises the probability of the task being completed, with the highest possible quality of service.

Human Derived Trust — To assess a trustee, a truster usually gathers evidence that supports one or more conclusions about the trustee’s likely behaviour. In the preceding examples, this evidence gathering can, at least in part, be automated. For instance, automatic intrusion detection indicates that a particular user account is being used for malicious purposes, or that a service provider may be judged on the quality of service it has provided in the past. In other cases, such as online auction houses like e-bay¹, trust may depend on intangible qualities, only discernible from the subjective experience of a human user.

1.3 Service-Oriented Computing

Although there are a number of areas in computing where trust is an issue, our main motivation is the recent interest in open systems that adopt a service-oriented architecture (Huhns and Singh, 2005). This type of architecture is central to the idea of web services (McIlraith et al., 2001; Roy and Ramanujan, 2001), pervasive computing (Adelstein et al., 2004), and grid computing (Foster and Kesselman, 2004). Each of these fields has arisen individually in response to different problem domains, but there is a great deal of overlap in the challenges they face. In particular, they involve loosely coupled software modules (or services), which are usually distributed over a network, and can be composed dynamically to solve different problems. These services may perform a variety of roles from implementing a specific algorithm, to mediating access to databases, storage devices or compute clusters.

In this section, we shall focus mostly on grid computing, because part of the work presented in this thesis is a system that targets this field in particular (see Section 4.5). However, the specific contributions that we claim in this thesis lay in the intersection

¹<http://www.ebay.com>

between grid computing and other service-oriented environments. Thus, we believe our work has general applicability within service-oriented computing.

From our perspective, what is interesting about all these systems is that resources belonging to different stakeholders can be brought together in support of a common goal, and that individual resources may become available, or unavailable, at different times. With regard to grid computing, Foster et al. (2001) state the aim of the Grid as facilitating, “*coordinated resource sharing and problem solving in dynamic, multi-institutional virtual organisations.*” Specifically, the Grid is concerned with direct access to computers, software, data and other resources for multiple purposes that involve collaboration across geographical and institutional boundaries.

In this context, a virtual organisation (VO) is the set of individuals or organisations that are involved in such a collaboration. It is envisaged that the resources available to a VO may offer varying degrees of reliability, and may leave and re-join the system at any time. In addition, the organisations that supply these resources could have different, and possibly conflicting, interests. Together, these properties imply an inherent unreliability in (individual) grid resources, which is amplified further when we consider that the Grid is intended to take on a global scale.

Despite the emphasis on dynamic behaviour within the Grid, most existing grid systems are somewhat inflexible in terms of their inter-operation and interactions. Instead, VOs have usually consisted of research institutes with longstanding relationships, whose computing resources have been used together in reasonably pre-defined ways. However, the long term goal is wider than this, and extends to the implementation of *ad hoc* solutions that bring together stake holders on-the-fly under short term relationships, potentially for financial gain.

Managing VOs that form and adapt rapidly is not something that has been adequately addressed by current Grid technology. The emphasis has instead been on developing protocols and middleware that allow secure, robust and scalable inter-operation between resources. Therefore, it has been suggested in Foster et al. (2004) that to fill this gap we should turn to autonomous agents and multi-agent systems (MAS).

As defined by Wooldridge and Jennings (1999), an agent is, “a self-contained problem-solving system capable of autonomous, reactive, proactive, social behaviour.” A multi-agent system is a system that comprises a number of such agents, which need to interact with each other to achieve their goals. Much like in grid computing, agents often need to organize themselves into collectives (or VOs), despite having potentially conflicting interests. However, the emphasis in agent research has been much less about building robust and scalable infrastructures, and much more on endowing agents with problem-solving abilities, so that they may interact and adapt in a changing environment, without human intervention.

These types of properties have led many authors to recognise the relevance of agent technologies to service-oriented computing as a whole. For example, there has been much research focused on agents that exchange goods in an open marketplace (Pardoe et al., 2006; He et al., 2006; Sierra, 2004), negotiate and participate in auctions with sophisticated strategies (Nguyen and Jennings, 2005; Gerding et al., 2006), and form resource-sharing coalitions of their own accord (Blankenburg et al., 2005; Dang and Jennings, 2006). As argued by Foster et al. (2004), these types of techniques make agents a strong candidate for resource management in grid computing, a role which is quite common in the literature with regard to other service-oriented environments (Huhns et al., 2005).

1.4 Research Objectives

It is clear from the previous section that multi-agent systems provide an appropriate tool for managing service-oriented computing environments, like the Grid. It should also be apparent that such systems may include many different stakeholders, who need to interact despite having potentially different incentives and goals. As such, issues of trust naturally arise from these systems, and decision makers need to be aware of them. If these decision makers are autonomous agents, then they should be endowed with the ability to reason about trust, so that they can make effective decisions.

It is this need that forms the motivation of this thesis and the basis of our work. In particular, we aim to *develop trust assessment models, or mechanisms, that can be used to aid decision making by autonomous agents in a service oriented environment*. In doing so, we can allow agents to manage resources in a service-oriented system effectively, by identifying which resources are most reliable, and which are best avoided. This information could then be used along with other factors, to allow an agent to choose the most suitable resources to fulfill its needs, and adapt to circumstances in which resources fail, or become available.

This raises two issues: (1) we should consider what kind of trust issues agents need to address, and (2) we should consider what factors an agent should account for, when making decisions involving trust.

To address the first of these issues, we refer back to the categories of trust issues outlined in Section 1.2. All three of these could potentially arise in problems that an agent may face, but not all will do so to the same extent. This is due to the nature of the problem solved by agents in a service-oriented environment. In general, we expect agents to interact directly with each other, rather than through users. Therefore, it is possible for agents to learn about each others' behaviour directly, and so the need to deal with human derived trust is perhaps of less importance. Similarly, security and trust management issues are of less interest to us, because they can generally be

addressed outside the realm of agents, and because security policies are still something that warrant significant human involvement. On the other hand, the key benefit of using agents in service-oriented environments is to manage the coordination and allocation of resources from different domains. This involves choosing which services to use, and for what purpose, so it is important to assess the relative reliability of potential service providers. Thus, the main type of trust that we are interested in is the ability to assess the willingness and capability of an agent to provide a particular service.

To address the second issue, we need to consider the information that is available to an agent when making decisions involving trust. This very much depends on the situation, but examples include using social rules and norms that apply in an environment (Ramchurn et al., 2003) and the relationships that are known to exist between agents (Ashri et al., 2005). Although these can play a significant role, perhaps the best indicator of how trustworthy an agent will be is how it has actually behaved in the past. Observations of past agent behaviour are thus widely recognised as an important and basic predictor in trust assessment (see Section 2.2 for a detailed review of the literature in this area).

That said, observations come in different guises, and not all are as reliable as each other. In particular, the source of information is important: a truster may observe a trustee's behaviour itself, or it may enquire about a third party's experiences with a trustee. Intuitively, direct experience is more reliable than second-hand opinions, but this may not always be available, especially if a system contains many agents who interact frequently with strangers.

From human society, we also know that third party experience can be a useful source of information, the sum of which we refer to as an agent's *reputation* (Sabater and Sierra, 2001). For example, we often make assessments about a new business based on word of mouth, or by reviews written in the press. This form of assessment also has a strong mediating effect on agent behaviour: if a business provides a poor service and it knows that this will affect its general reputation, then this will increase its incentive to provide a good service to continue receiving custom in the future.

On the other hand, those that provide opinions may not always have a truster's best intentions at heart, or may not evaluate a service in the way that a truster expects. For example, if a reputation source is asked to rate a trustee's performance in an area in which the reputation source competes, it may provide an unfairly critical report, so as to increase its own market share. In addition, an agent's environment may mean that its performance is perceived differently depending on which agent it supplies a service to. In this case, a trustee's performance, as observed by a reputation source, may not be a good indicator of its performance toward a particular truster. As such, reputation cannot be considered as reliable as a truster's direct experience, and so it must guard against the possibility that reputation is inaccurate, when making its decisions.

With this in mind, we can now set out a number of objectives, which we intend to fulfill so that we may meet the aim stated at the start of this section. Specifically, we believe the following six objectives are important.

1. **Decision Facilitation** As our main aim is to facilitate an agent's decisions, there should be a clear mechanism by which any system we develop can be used to this end. Although this seems obvious, it is conceivable to construct a system that measures some concept of trust, for which the precise role in decision making is not well defined.
2. **Assessment Based on Direct Experience** In very many cases, it is safe to assume that a trustee's past behaviour is a good indicator of its future behaviour. Therefore, if a truster has previous direct experience with a trustee, it should account for this information when deciding how to interact with it in the future. As such, we should facilitate this in our work.
3. **Assessment Based on Reputation** In the absence of direct experience, third party experience can be a useful indicator of a trustee's behaviour. Thus, our trust model should make use of a trustee's reputation, based on third party experiences.
4. **Robustness against Inaccurate Reputation Sources** Sources of reputation will not always be reliable. For example, a close colleague of a trustee is likely to have a strong incentive to exaggerate the trustee's credentials, and hence provide unreliable information. Therefore, while making use of reputation, our trust model should be robust against inaccurate information sources.
5. **Efficient Opinion Communication** To use reputation, third parties must communicate their knowledge to a requesting truster. Hence, we need to support this communication, ideally such that all relevant information is transmitted with minimum cost.
6. **Assessment with Varying Degrees of Evidence** As a group of entities interacts over time, the number and type of interactions that occur between group members may change. Our trust assessment model should make use of this information, on average increasing the accuracy of its results as the frequency of interactions in the system increases. However, the system should not be dependent on such information, but instead its performance should degrade gracefully as information decreases. In particular, the system should be able to operate when any one of the following statements is correct:
 - The truster has direct experience of the trustee.
 - The trustee is not known directly by the truster, but is known by other entities within the system.
 - The trustee is not known to any entity in the system.

In each case, we should endeavour to make as much use of the information provided by the environment as possible, but should also be able to provide reasonable results when certain sources are not available. This is particularly important, for instance, when a large system has just been initialised and no interactions have taken place.

Together, these objectives constitute a set of requirements that a trust model must achieve, if it is to facilitate effective decision making in the types of environments outlined in the previous section. As should become clear from Chapter 2, each of these requirements is fulfilled, at least in part, by mechanisms described in the existing literature. However, these existing systems all suffer from certain limitations, which are addressed in this thesis. Specifically, we make a number of contributions to the state of the art, which we outline in the next section.

1.5 Research Contributions

One factor that determines how the objectives defined in the previous section can be addressed, is how the actions of a trustee can be appropriately represented. This depends on the needs and preferences of agents within the target domain. For example, suppose you ask someone to fetch a pint of milk from the shop before 8.00pm. In this case, it may be that you only care that the milk is delivered by that time, in which case the relevant aspect of behaviour can be represented as a binary event — either the milk is delivered, or it is not. On the other hand, you may care about more graded aspects of this service, for example, the price of the milk, or the number of days left before it becomes unfit for consumption.

With this in mind, the main contribution of our work lies in the creation of two computational models of trust, which each address our objectives for different representations of trustee behaviour. In particular, the first of these, known as TRAVOS (Trust and Reputation system for Agent Based Virtual OrganisationS), is designed specifically for cases in which trustee behaviour can be represented as a binary event. In contrast, the second model, known as TRAVOS-C, extends the basic concepts employed by TRAVOS, to address cases in which behaviour is represented as a continuous set of real numbers².

While together these models do not cover the complete set of possible behaviour representations, these two cases do constitute an important subset, which deserve separate treatment. In addition, the lessons learned from both models provide insight into how trust could be addressed in other domains. In the rest of the section, we detail the main contributions that these models share, and exhibit individually.

²The C in the name refers to the model's applicability to continuous action spaces.

What is common to both models is that they both use probability theory to represent and derive values of trust, regarding the future action of an agent. As such, they share with other probabilistic models two key advantages over trust models that use alternative representations. First, the known rules and results of probability theory can be used to derive a number of important properties of these models. For example, estimates of unknown quantities about an agent's behaviour can be shown to be optimal according to certain criteria and under certain assumptions. Second, probability values have a natural interpretation in decision theory, so can be used in a well defined and understood manner to guide rational decision making. However, in addition to these standard benefits, we provide the following six advances over the state-of-the-art with respect to other probabilistic models of trust.

1. We specify the first set of general requirements for communicating reputation that, if satisfied, guarantee three properties: (1) if a truster requests information about a trustee from an honest third party, then the truster will receive all relevant information about that agent's experiences; (2) this is done such that, under certain conditions, estimates based on reputation are consistent and as reliable as estimates based on equivalent direct experience; and (3) this is achieved with minimal communication overhead. Moreover, this set of requirements is shown to be satisfied by both of our models, TRAVOS and TRAVOS-C.
2. In TRAVOS, we specify a model of trust for assessing trustees based on a binary representation of behaviour. This provides a mechanism for filtering out inaccurate reputation that, of all previous models of its kind, is shown to outperform the only existing method for achieving this. A description of the TRAVOS model, along with an explanation of its usage, is given by [Patel et al. \(2005a, 2004\)](#), while empirical results comparing the model's filtering mechanism to its existing competitor is given by [Teacy et al. \(2005, 2006\)](#).
3. Using TRAVOS, we are the first to show how a trust model can be used, as part of a multi-agent system, to manage the formation and reformation of virtual organisations, within a grid computing environment. To do this, we deployed TRAVOS as part of a grid resource management system, known as CONOISE-G (Constraint Oriented Negotiation in Open Information Seeking Environments for the Grid). Details of this system, including the role of trust within the system, are given by [Shao et al. \(2004\)](#); [Patel et al. \(2005b,c, 2006\)](#); [Nguyen et al. \(2006\)](#).
4. In TRAVOS-C, we specify the first model for assessing trustees based on a continuous representation of behaviour, which satisfies the communication requirement, outlined in Point 1.
5. As part of TRAVOS-C, we have produced the first method for dealing with inaccurate reputation, which is derived completely using probability theory. As a direct result, this method is known to be optimal under the model's assumptions.

6. It has been suggested by some that, when a truster has little or no direct or indirect experience of a trustee, predictions about the trustee can be improved by taking into account the known behaviour of similar agents in the system. We have developed such a mechanism in TRAVOS-C, which is the first such method for a probabilistic trust model with continuous action spaces. Moreover, this is the only mechanism of its kind which automatically adapts its impact on prediction to best suit conditions in the truster's environment.

1.6 Thesis Structure

In the rest of the thesis, we survey the existing literature on trust in multi-agent systems, and describe in detail the work through which we have made the contributions outlined in the previous section. This is achieved through the course of the remaining chapters, which are structured as follows:

- Chapter 2 gives a review of existing mechanisms for managing trust in multi-agent systems, concentrating in particular on work that goes some way towards fulfilling the objectives described in Section 1.4. The most relevant models with regard to these objectives are identified, and their prominent characteristics are described. In light of these characteristics, we identify several significant limitations of existing trust models, and use these to motivate the work described in the proceeding chapters.
- Chapter 3 defines a framework for reasoning about trust, and defines the basic notation used throughout the rest of the thesis. Through this framework, we show how, in general terms, a group of agents can make rational decisions regarding their peers, based on observations of past behaviour. Moreover, in cases where these observations are shared between agents, the framework shows how this can be achieved, while fulfilling our objective of efficient communication, outlined in Section 1.4.
- Chapter 4 introduces the TRAVOS system, describing how it instantiates the framework from Chapter 3 to reason about trustees whose behaviour is represented as a binary event. Furthermore, we show how TRAVOS can be used in practice, by describing how it operates within the CONOISE-G system, and by giving an empirical evaluation of its performance. In particular, this evaluation demonstrates how the reputation filtering mechanism in TRAVOS outperforms its nearest competitor.
- Chapter 5 introduces TRAVOS-C, describing how it too realizes the framework from Chapter 3 by building on the characteristics of TRAVOS. In particular, we describe how it can be used to reason about trustees with continuous behaviour

representations, outline its advantages over both TRAVOS and other existing trust models, and detail its theoretical properties. Finally, we give an empirical evaluation of TRAVOS-C, demonstrating how it performs under the model's assumptions, and how it is robust against certain violations of those assumptions.

- Finally, Chapter 6 summarizes the conclusions drawn throughout the thesis, and in particular the contributions and limitations of the techniques we have developed. In addition, we describe and motivate our plans for future work.

Chapter 2

Computational Models of Trust

As stated previously, there are a number of areas for which trust is relevant. Generally, these fall into three major categories: security issues, human derived trust and service provision. In our work, we are concerned primarily with service provision, and so in this section we focus our attention only on related work that is relevant to this set of problems. Specifically, we divide our discussion into five main sections. First, Section 2.1 considers mechanisms for assessing trust by analysing the different types of beliefs an agent must hold to have trust in another. Second, Section 2.2 surveys the major approaches for forming trust based on information directly available to a truster. Third, Section 2.3 addresses the problems that arise when trust is based on third party opinions. Fourth, Section 2.4 outlines an alternative approach to trust assessment: it discusses mechanisms that attempt to enforce trustworthy behaviour by making it in a trustee's best interest to be trustworthy. Finally, we conclude the chapter in Section 2.6, summarising the related work to date, and identifying key outstanding issues.

2.1 The Cognitive Viewpoint

One way of analysing the trust that one agent should have in another, is to reason about the beliefs that a truster should hold about a trustee, if it is to believe that the trustee will behave in a particular way. In this respect, we consider the work of [Castelfranchi and Falcone \(2001\)](#), who adopt the same basic definition of trust as us (i.e. a subjective probability of a trustee's action). Building upon this, they make two things explicit: (1) they identify beliefs that a truster must hold before it can rationally believe a trustee will carry out a given action; and (2) they identify a three-way relationship that exists between the concepts of *trust*, *fear* and *authority*. The core beliefs that are prerequisite to a belief of trust are as follows:

- The truster must believe that the trustee is *willing* to carry out the action.

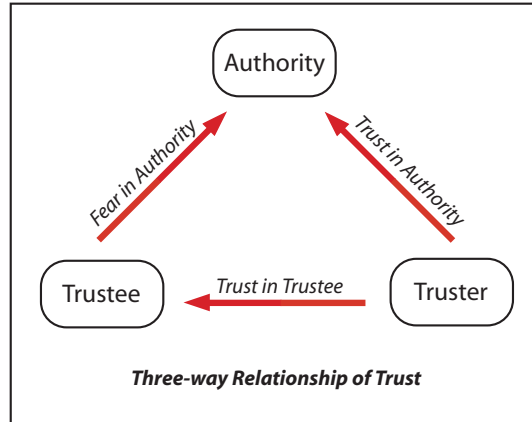


FIGURE 2.1: The three way relationship between trust, fear and authority.

- The truster must believe that the trustee is *capable* of carrying out the action.

In turn, these beliefs may be conditioned on a number of other beliefs that, for the most part, are domain dependent. In general, however, we can distinguish between two different sets of beliefs: internal beliefs, which relate to the trustee’s mental state, and external beliefs, which concern environmental conditions. To illustrate the impact of the latter, consider an entity, Captain Joe, who is capable of sailing a particular boat, the Jolly Roger. If something was to happen to the Jolly Roger so as to cause it to sink, then Captain Joe will no longer be able to sail the boat, despite his skills as a sailor.

Figure 2.1 illustrates the relationship between trust, fear and authority mentioned earlier. Fear can be said to be negative trust; it is trust in an entity to carry out an action that has a negative effect on the truster. Like trust, fear in an entity requires the conditions of willingness and capability to be present. When we introduce an authority, which is capable of punishing unsolicited behaviour, an interesting dynamic is set up between the authority, truster, and trustee. The fear a potential wrong-doer has in an authority decreases its likelihood of behaving illegally. On the other hand, if a victim trusts an authority to protect its rights, and it can assume potential criminals hold similar beliefs about the authority, then its trust in potential criminals can be increased.

The influence of authority on trust relationships is also acknowledged by Dasgupta (1988), who argues that, if a rational agent is put in a position where it can choose to benefit at the expense of others, it will always choose to benefit unless it has reason to fear retribution.

With these factors identified, we need to consider how can we use them to develop an automated method for reasoning about trust. An attempt to do this is made by Falcone et al. (2003), in which they use fuzzy logic (Zadeh, 1975, 1965) to build a truster’s beliefs about trust, willingness and capability from other beliefs that are largely domain dependent.

Although identifying the composite beliefs that make up trust gives us a better understanding of what we are attempting to measure, one major question remains unanswered: how can a trustor determine its core beliefs based on observations of its environment? Clearly, factors such as trustee willingness and capability are not directly observable in general; they must be estimated from observable evidence. Moreover, depending on evidence that is observable, attempting to estimate separate beliefs about such factors may not be practical at all. For example, consider a scenario in which all we can observe is an agent's external behaviour in the absence of any other environmental data. In this case, the best we can do is quantify the uncertainty in the trustee's behaviour directly; we cannot possibly distinguish between the trustee defaulting on its obligations because it wants to, or because it cannot do otherwise.

Thus, due to the difficulty in obtaining information about an agent's internal beliefs, in our work, we choose not to employ techniques for reasoning about such beliefs, but to rely only on evidence pertaining to an agent's external behaviour.

2.2 Learning from Direct Observations

In this section we turn our attention to methods of representing trust, and how to ground such representations in evidence directly observable to a trustor. We differentiate direct evidence from evidence as reported by other agents, the latter of which raises a separate group of problems that we address in Section 2.3.

Generally, existing trust models represent trust in one of three ways: (1) they adopt an improvised representation, based on intuitive assumptions about the meaning of trust; (2) they apply probability theory; or (3) they apply Dempster-Shafer theory. For the purposes of clarity, we separate our discussion according to this categorisation, and discuss each in turn in the subsequent subsections.

2.2.1 Improved Models of Trust

As mentioned earlier, although the concept of trust is prevalent in society, there is some disagreement and confusion about its precise definition. Perhaps partially as a result, a range of different representations have been adopted in existing computational models of trust. In some cases, trust is modelled as belonging to a finite set of qualitative labels, examples of which include the work by [Azzedin and Maheswaran \(2002c,b,a\)](#) and [Abdul-Rahman and Hailes \(1997\)](#). In the case of the former, the trust of one entity in another is a value belonging to the set $\{A, B, C, D, E, F\}$, and similarly in the latter, a member of the set $\{-1, 0, 1, 2, 3, 4\}$. Typically, these values are associated with linguistic labels that describe their intended meaning. For instance, Abdul-Rahman and Hailes attach labels to trust values as described in Table 2.1.

<i>Value</i>	<i>Meaning</i>	<i>Description</i>
-1	Distrust	Completely untrustworthy
0	Ignorance	Cannot make trust-related judgement about entity
1	Minimal	Lowest possible trust
2	Average	Moderate trustworthiness
3	Good	More trustworthy than most entities
4	Complete	Completely trust this entity

TABLE 2.1: Trust value semantics used by Abdul-Rahman *et al.*

This relatively coarse set of values reflects a perceived difficulty in choosing continuous trust values with any meaningful degree of accuracy. In our view, however, this problem is specific to cases in which trust is elicited from a user¹. As should become clear from what follows, there are meaningful methods of calculating continuous trust values when trust assessment becomes a fully automated task. We therefore argue that the difficulty in distinguishing between discrete trust levels compared to continuous levels limits the former’s applicability to human elicited trust values.

Models that represent trust as a real-valued scalar include those developed by Marsh (1994) and Zacharia *et al.* (1999); and more recently by Fan *et al.* (2005), and Griffiths and Chao (2005). Representative of these, and one that makes a good attempt to deal with the objectives outlined in Section 1.4, is the REGRET system (Sabater and Sierra, 2001, 2002), which includes three dimensions of trust: an individual dimension, a social dimension, and an ontological dimension. We shall examine each of these in turn below.

First, we consider the ontological dimension that is essentially concerned with the subjectivity of trust with respect to an individual truster. A trust value in REGRET is represented as a numeric value in the range $[-1, 1]$, with a value of 1 interpreted as *absolute* trust, and -1 interpreted as complete distrust. These values are attached to a particular context by a label, examples of which are `to_overcharge`, meaning that a trustee has a tendency to charge more for a service than the truster believes it is worth, and `quality_swindler`, meaning that, from the perspective of the truster, the trustee tends to supply services with unacceptable quality. The intended interpretation is that they relate to a particular trait of a trustee’s behaviour.

An important element of the ontological dimension is that behavioural traits² of an agent can be defined in terms of other, lower level traits. For instance, a service provider could be assessed according to a trait labelled `swindler`, which is defined in terms of the traits, `to_overcharge` and `quality_swindler`, mentioned earlier. From an implementation perspective, REGRET calculates trust values for such compositional traits as a weighted average of the trust values calculated for the base traits. The weights used in this

¹To illustrate, consider trying to assess the probability of it raining tomorrow; is it possible to decide whether this probability is more likely to be 2.1 or 2.2?

²In REGRET, these are referred to as ‘reputation types’. We use the term trait, so as not to confuse it with the concept of reputation as a collective opinion of a group with regard to a particular entity.

calculation are considered to be dependent on an individual truster; they encode the agent's subjective definitions for these terms. Besides specifying how compositional traits can have trust values calculated, the ontological dimension can serve a communication role in that a reputation source can share its definition of compositional traits with other agents, so that they can decide how best to interpret reputation information from that provider.

The individual dimension of trust is based solely on the first-hand knowledge that a truster has about a trustee. In REGRET this is calculated based on past interactions that have occurred between the truster and trustee. For example, when a truster purchases a service from a trustee, the truster will have expectations about how the trustee will behave. Some of these expectations will be explicit, based on a contract between the truster and trustee for what the trustee should provide. Others will be implicit, based on the trustee's personal perspective on the world. A truster's individual trust level (with respect to a particular trait) is a function of the difference between the utility the truster would achieve if the trustee behaved according to these expectations and the actual utility gained from the interaction.

As well as providing a method for calculating these trust values, REGRET also provides two separate methods to measure the *reliability* of these values. Two different types of uncertainty determine the reliability of a trust value:

- Intrinsic uncertainty in trustee behaviour, which is estimated based on the variance of observed trustee behaviour.
- Uncertainty due to lack of evidence, for which REGRET uses a function that decreases logarithmically until a minimum value, in line with the number of observed interactions with a trustee and the time that has passed since those interactions.

In REGRET, both of these are measured using improvised functions. For example, REGRET calculates evidential uncertainty using Equation 2.1, adapted from Sabater and Sierra (2001).

$$\text{evidential uncertainty} = \begin{cases} \sin\left(\frac{\pi}{2 \cdot itm} \cdot noObs\right) & noObs \in [0, itm] \\ 1 & \text{otherwise} \end{cases} \quad (2.1)$$

Here, *noObs* is the number of observations a truster has made of a trustee's behaviour, and *itm* is a threshold number of observations, above which the truster considers its knowledge of a trustee to be completely reliable. In a similar way, intrinsic uncertainty is measured by another function, improvised from intuitive conclusions about what factors should affect its value. An overall reliability factor is then calculated as a weighted average of these two functions. One problem with this scheme is that it is unclear what value should be chosen for *itm*, and what weight should be used to generate the overall uncertainty value.

Often, an agent will need to assess its trust in an entity with which it has little or no previous experience. In this case, REGRET can draw upon the social dimension of trust, and there are three sources of information that fall under this heading: witness reputation, neighbourhood reputation, and system reputation.

Witness reputation is, as the name suggests, based on the opinions of third parties concerning a trustee. The influence of particular witness's opinion on the overall trust value is partly determined by the truster's trust in the reputation source to provide reliable information. This can be calculated by applying the formulae for individual trust, effectively treating the ability to give reliable reputation information as just another trait.

Neighbourhood reputation assumes that the truster maintains a *sociogram*, which is a network structure describing the social relationships between agents in the environment. To calculate neighbourhood reputation, REGRET applies a set of static³ fuzzy rules, where the antecedent of each rule is a condition on the relationships connecting the trustee to other agents. To illustrate, we might define a rule such as

$$\text{IF } \textit{coop}(\textit{trustee}, \textit{agent}.b) = \textit{high} \text{ THEN } \textit{socialTrust} = \textit{very_bad},$$

where *high*, and *very_bad* are predefined fuzzy sets.

System reputation is also calculated according to a static set of fuzzy rules. In this case, the rules are defined according to the role that a trustee plays within an institutional structure — *seller* is an example of such a role. As with neighbourhood trust, system trust assumes that information about social roles is available to the truster.

REGRET combines these different sources — individual and social dimensions — based on reliability functions defined over them. In addition, there is also an intrinsic preference ordering built in: direct interactions are intrinsically more reliable than witness reputation, witness reputation is more reliable than neighbourhood reputation, and neighbourhood reputation is more reliable than system reputation.

From our perspective, REGRET is significant because it broadly satisfies the objectives specified in Section 1.4. However, REGRET suffers limitations for at least two reasons. First, it assumes that certain information (for example, a sociogram) is available. Second, there are several parameters in the model, optimal values for which are not known, and may be domain dependent.

³By static, we mean that REGRET must be preconfigured with a set of rules. REGRET cannot learn these rules for itself.

2.2.2 Probabilistic Models of Trust

Aside from fuzzy logic, the trust models we have looked at so far all make use of, essentially, hand-crafted representations of trust, and operations defined on these representations. This is by no means an invalid approach — ultimately, the goal of assessing trust is to provide discriminatory information about trustees, and so any technique that can be shown to do this can be considered reasonable. However, there are existing formalisms for reasoning about uncertainty, which have well known beneficial properties, and are well grounded in mathematical theory. Of these, perhaps the most prominent is probability theory.

One noteworthy probabilistic trust model is detailed by [Barber and Kim \(2001\)](#). It provides a well grounded method for assessing the reliability of information sources, and shows how it can be used to combine conflicting information into a consistent knowledge base. Unfortunately, the model is designed specifically to deal with such conflicts: it uses the statistical properties of the conflicts themselves to perform its task, and so cannot be applied to a more general setting.

More generally, the majority of probabilistic models that attempt to assess trustees on a broad range of services (which include those reviewed in the remainder of this section) have two things in common. First, they represent the outcome of an interaction with a trustee as a bistable event — either the trustee cooperates and fulfills its obligations to the truster, or it defects and does not. Second, they estimate the probability distribution for this binary variable based on direct or indirect (via reputation) observations of the trustee's past behaviour. Obviously, this simplification leaves clear room for improvement: if a truster's utility is dependent, not only on whether a trustee performs a task, but also on how well the task is performed, then a bistable representation will fail to capture the relevant dynamics of the problem. Nevertheless, situations in which task performance does not carry much significance over and above task completion constitute an important subcase.

An example of such a system can be found in [Wang and Vassileva \(2003\)](#). Here, a trust mechanism is presented for use in a peer-to-peer file sharing environment. Trust in a particular provider is assessed according to several quality attributes, such as type of file requested, download speed and file quality. The system uses a naïve Bayesian network, in which the probability of the provider being trustworthy (modelled as a binary variable) is dependent on each of the quality attributes considered. Here, *naïve* means that the effect of each attribute on the trustworthiness of a provider is assumed to be independent. Such assumptions are often made to simplify a problem domain, with solutions adopting them generally being robust when faced with minor violations. Whether the assumption is reasonable in the domain targeted by this model depends on the particular set of quality attributes used in a given instance of the model.

	Agent <i>A</i>	Agent <i>B</i>
successful	20	2
unsuccessful	20	2

TABLE 2.2: Frequencies of successful and unsuccessful interactions with different agents.

One factor that Wang and Vassileva fail to account for is *evidential uncertainty*. Here, we differentiate evidential uncertainty from intrinsic uncertainty. We define intrinsic uncertainty to be uncertainty that is due to inherent unpredictability of a stochastic process. On the other hand, we consider evidential uncertainty as uncertainty that is due to a lack of knowledge. To illustrate, consider observing successful and unsuccessful interactions with two agents, *A* and *B*, the frequencies for which are shown in Table 2.2. Using Wang’s model, we would consider there to be no difference in the uncertainty surrounding agent *A*’s behaviour and agent *B*’s behaviour. However, intuition tells us that this is not the case, because we have interacted with *A* ten times more than *B* and can therefore be more certain about *A*’s true behaviour. This highlights a failing common to all simple probabilities that is particularly important in domains where the frequency of observations is relatively low. We believe that trust assessment in large multi-agent systems is such a domain, because the likelihood of any two agents interacting a large number of times is fairly low. We therefore argue that accounting for both types of uncertainty is important and give further justification for this in Chapter 3.

Fortunately, to say that probability theory in general cannot account for evidential uncertainty would be incorrect. This is illustrated by the trust model presented by Jøsang and Ismail (2002), in which trust is modelled as a probability distribution for a binary event, a class of distributions commonly referred to as *Bernoulli* distributions (Evans et al., 2000). In addition, however, they also model the *parameter* distribution of the Bernoulli distribution (DeGroot and Schervish, 2002). Statistical models, such as Bernoulli distributions, are characterised by a set of parameters that determine their shape. In the case of a Bernoulli distribution, it is characterised by a single parameter — the probability of the variable being equal to one. The parameter distribution in this case is the distribution over the possible values of that probability.

For simplicity, the authors choose to represent the parameter distribution as a beta distribution. The advantage of this is that there is a special relationship between Bernoulli and beta distributions. Specifically, consider a Bernoulli distribution for which the prior parameter distribution is a beta distribution. If we draw samples from this Bernoulli distribution under an i.i.d assumption⁴, then the posterior parameter distribution, given the samples, will also be a beta distribution. A family of distributions which exhibits this property for a statistical model is known as the model’s *conjugate* family.

⁴This is a standard abbreviation for the assumption that samples are drawn independently from an identical distribution.

Effectively, the parameter distribution represents the evidential uncertainty surrounding the true intrinsic probability of a random variable. In this case, the intrinsic probability that a trustee will cooperate rather than defect. It can be used to reason about how much evidence is required to make a particular decision, or to compare the confidence levels different agents have in their knowledge about a trustee. Again, we discuss this further in Chapter 3. Moreover, by choosing a conjugate prior, the authors simplify the process of calculating, combining, and storing the parameter distribution associated with a trustee. For this reason, beta distributions are also applied to the field of trust by Klos and Poutré (2004); Mui et al. (2001); Zhang and Cohen (2006) and Buchegger and Boudec (2003).

2.2.3 Dempster-Shafer Models of Trust

An alternative method for handling uncertainty can be found in Dempster-Shafer theory (Shafer, 1976). Dempster-Shafer provides a mechanism for forming degrees of belief about sets of hypotheses, based on available evidence. For example, imagine we have a set of two competing hypotheses $\{A, B\}$, of which only one can be true. Dempster-Shafer theory divides the total belief in the set between the elements of its superset⁵, $\{\{A\}, \{B\}, \{A, B\}\}$. Essentially, belief in the set $\{A\}$ represents the evidence supporting A as the true hypothesis (and likewise for set $\{B\}$). On the other hand, belief in the set $\{A, B\}$ is belief that cannot be divided between A and B . This can be said to represent the evidential uncertainty surrounding A and B ; because of this, Dempster-Shafer theory is often claimed as a solution to the inability of simple probabilities to capture this notion.

In particular, this is the rationale given for its use by Yu and Singh (2002). Here, the authors define a binary hypothesis set, in which the competing hypotheses are that an agent is trustworthy, and that it is not trustworthy. They consider scenarios in which trustees supply services, which are given a *quality of service* rating between 0 and 1. To gather evidence for the trustworthiness of an agent, they perform the following three steps. First, they break the range of quality values into three intervals, $[0 \leq x < a]$, $[a \leq x < b]$, $[b \leq x \leq 1]$, where a and b are arbitrary constants. Second, they count the proportion of recent⁶ trustee interaction outcomes that fall in each of these three intervals. Finally, they take the proportion of outcomes in the lower interval as the belief that the trustee is untrustworthy, the proportion in the higher interval as the belief that the truster is trustworthy, and the proportion in the middle interval as the belief in the total set, $\{\text{untrustworthy}, \text{trustworthy}\}$. In line with

⁵By definition, the superset of a set, S , is the set comprised of all subsets of S .

⁶In their model, Yu and Singh only use the x most recent observations of a trustee's behaviour. This allows for the possibility that a trustee's behaviour changes over time, in which case old observations would be poor predictors of behaviour.

Dempster-Shafer theory, belief in the total set is interpreted as the degree of uncertainty in whether the trustee is trustworthy or not.

The problem with this approach is twofold. First, there is no clear way to choose the values for the constants a and b . Second, the notion of trust that this representation captures is somewhat artificial. Consider as an example a trustee with whom a truster has (recently) interacted 1000 times. On each of these occasions, the trustee's quality of service was precisely 0.5. Here, the truster has chosen $a = 0.3$ and $b = 0.7$. According to Yu and Singh's model, this means that the truster is completely uncertain whether it trusts the trustee or not. A more useful conclusion would be that the trustee provides an average quality of service of 0.5, with very low variance, so that there is a large degree of certainty regarding its behaviour.

Furthermore, there are three reasons why Yu and Singh's rationale for using Dempster-Shafer theory, rather than probability theory, is somewhat unsound. First, despite their claim to the contrary, probability theory can be used to represent and reason about evidential uncertainty, through the use of parameter distributions (Section 2.2.2).

Second, it is well known that Dempster-Shafer theory can result in counter-intuitive conclusions if there is a significant degree of conflicting evidence surrounding possible world states (Campos and Cavalcante, 2003). With regard to trust assessment, this is a significant problem because it is highly possible that a truster may receive conflicting opinions about a trustee, particularly if those opinions come from agents with conflicting interests that may provide misleading evidence in support of their own goals.

Finally, Dempster-Shafer theory is best suited to applications in which available evidence supports sets of hypotheses to an equal degree, rather than individual conclusions. If this is not the case, then probability theory can provide a more simple and intuitive method of reasoning. For example, if an agent provides a QoS of 0.7 this supports the hypothesis that it will provide a QoS of 0.7 in the future, and to a lesser degree QoS values close to 0.7. Conversely, there is no QoS value that a trustee could provide that would equally support the hypotheses that a trustee will provide QoS values of both 0.1 and 0.9, which supports the use of probability theory rather than Dempster-Shafer theory.

An alternative application of Dempster-Shafer, relevant to trust, is given by Jøsang (2002, 2001). Here, it is used to define *subjective logic*, which is an attempt to extend first order logic to reason about propositions that have probabilities attached to their truth or falsehood. This has grounding in both probability theory and Dempster-Shafer theory, and has first order logic as a special case. Significantly, from the perspective of trust, it defines two new operators for reasoning about third party opinions: the *consensus* operator and the *discounting* operator. In particular, these operators can be used to combine opinions from different sources, as is required when trust is based on the opinions of others (Section 2.3). The consensus operator is used to combine opinions

from different sources when each source is equally and fully trusted to provide accurate information. The discounting operator plays a supportive role to the consensus operator: it is applied prior to the consensus operator, to any sources which are not fully trusted to provide accurate information, and its effect is to increase the evidential uncertainty surrounding the opinion. As a result it decreases the effect it would otherwise have when combined with other opinions.

The justification for these two operators, is grounded in statistical theory. Specifically, a mapping is provided between the Dempster-Shafer notion of evidential uncertainty, and the beta distribution representation described in Section 2.2.2. The operators are thus shown to be equivalent to operations on the parameter distribution. In the case of the consensus operator, the grounding relation first assumes that each opinion concerns a Bernoulli distribution, and that they are each based on disjoint sets of samples from that distribution, under an i.i.d assumption. The result of the consensus operator is then shown to be equivalent to the probability that would result, if all the data are considered together. Although the assumptions behind this grounding are not expected to hold in general, it is expected to give reasonable results, even when they do not hold. The discounting operator is given a similar justification, which we do not describe in full here. Briefly, under certain conditions, it is shown to be equivalent to multiplying an opinion by the probability that it is true.

Overall, subjective logic provides a promising method for reasoning about uncertain probabilities. In particular, its grounding in probability theory gives a good justification for its use. There are, however, two points that must be considered. First, the discounting operator does not say how the probability that a source is accurate should be calculated. This is an open question that may depend on the type of information available. Second, the definition is subjective, particularly in the case of its consensus and discounting operators, which make certain assumptions that may not be appropriate in every case. These should be questioned with respect to any application for which subjective logic is considered.

2.3 Learning from Others

The basic problem of trust assessment is to estimate the behaviour of a trustee based on the available evidence. When this evidence is gathered indirectly via third party opinions (reputation), there are four additional factors that we must consider. First, a third party may define observed properties in a different way from the truster. For instance, what one agent considers a good service may not be what another considers good. Second, reports from several different reputation sources may be based on the same evidence, resulting in correlated evidence. Third, the behaviour of a trustee towards a third party may be different from its behaviour towards the truster. Fourth, a

reputation source may have no incentive to provide reputation or, if it does, it may have an incentive to misrepresent its knowledge about a trustee. We can subdivide the latter into positive discrimination (collusion) in which the reputation source overestimates the beneficial qualities of a trustee, and negative discrimination in which the trustee's beneficial qualities are underestimated.

Each of these factors manifests itself as a decrease in the predictive power of reputation when compared to direct evidence. Thus, many trust models which employ reputation include bias reduction mechanisms to target one or more of these factors. Essentially, there are three basic methods for doing this, which we review separately in the subsections that follow. Specifically, Section 2.3.1 examines models that employ external factors, which are not directly related to a reputation source's opinions, to assess the accuracy of those opinions; Section 2.3.2 considers methods that assume that the majority of opinions received about a trustee are reliable; and finally, Section 2.3.3 reviews methods that assess the perceived accuracy of past reports given subsequent trustee behaviour.

2.3.1 External Factors

Zacharia et al. (1999) introduce two complementary reputation systems, HISTOS and SPORAS. Of these, SPORAS is a simple trust model that is not context dependent, and that a truster can use when there is little information available about other agents. To account for the unreliability of a reputation source, it simply weights its opinion by the truster's trust in the reputation source itself. The implicit assumption here is that if an agent can be trusted in general — for example, to provide a particular service — then it can be trusted to provide accurate information about other agents. Clearly, this assumption does not hold in general.

HISTOS, on the other hand, is a more sophisticated model suited to environments in which more information about a trustee's peers are available. It suffers from the same context independence as SPORAS, but takes on board the social relationships that exist between the truster, its reputations sources, and the trustee. Specifically, it (like SPORAS) treats trust as a transitive relationship, in which the trust of a truster in a trustee is a function of the trust of each reputation source in the trustee, and the trust of the truster in each reputation source. Unlike SPORAS, HISTOS builds a social network from the pairwise ratings that have previously been reported between agents. This is a directed graph, in which agents are represented by nodes, and edges between nodes represent the direct trust value of the parent node in the child node. The transitive trust relationship is then applied recursively along the directed paths between truster and trustee.

The REGRET system (Section 2.2.1) applies two different techniques to reputation noise reduction. The first of these applies the same transitive reasoning to trust as HISTOS and SPORAS, weighting a reputation source's opinion by the trust the truster has in that reputation source. However, REGRET's notion of trust is more expressive: it takes on board contextual issues such as the time a rating was given, and through its ontological dimension, can account for several different aspects of trust. For example, the trust of an agent as a reputation source may be built on its trust as a service provider and the accuracy of any past opinions it has provided. Unfortunately, REGRET does not give specific guidance on the relative importance of such factors, nor how the accuracy of past opinions should be calculated, the latter of which, in itself, is not a trivial issue.

The second mechanism adopted by REGRET specifically attempts to deal with the correlated evidence problem, which the majority of trust models recognise as a potential problem. The universal solution is to specify that a reputation source should only share its direct knowledge, and not pass on other agents' knowledge as its own, but the specifics of the solution vary. Most models assume independence, which can be justified if the intersection between agents' world views are small. In contrast, REGRET's solution is to carefully select reputation sources based on social network analysis. To do this, it uses an algorithm that divides a social network into groups of agents and then chooses reputation sources which are representative of those groups. The intuition is that a highly connected group of agents are likely to share the same knowledge, whereas loosely connected individuals are unlikely to share knowledge.

Common to all of these solutions is that they suffer from at least one of two problems: they assume that the information they rely on is readily available, which is not necessarily the case; or the links made between the factors they use and reputation accuracy do not hold in many cases. Specifically, information about the relationships between agents may not be readily available, or may require significant user input; and trust in an agent in one context (such as service provision) should not imply trust to provide reliable opinions. To rectify these limitations, the techniques described in the following subsections attempt to estimate accuracy based on the observed behaviour of both reputation sources and other agents.

2.3.2 The Majority Opinion

One way to judge opinion accuracy is to assume that the majority of opinions received about a trustee are reliable, and so discredit reports that deviate significantly from mainstream opinion. Such methods have been described by Jøsang et al. (2005) and Whitby et al. (2004) as *endogenous*, because they rely only on the statistical properties of the opinions themselves, while approaches that do rely on other factors are described as *exogenous*. Representative endogenous techniques are given by Whitby et al. (2004), and Dellarocas (2000), which we describe below.

Whitby *et al* extend the Beta Reputation System (Section 2.2.2) by applying an iterative filtering algorithm. In each cycle, an interquantile range⁷ is calculated for the set of opinions received about a trustee. Any opinions lying outside this range — that is, opinions that deviate significantly from the mean — are discarded. In the following cycle, the interquantile range is recalculated without the discarded opinions, and the process continues until all remaining opinions are in range. Although this approach is reasonable, and has been shown to give encouraging results, there is no guarantee that any opinions will remain after the algorithm has been applied. This can occur when all opinions differ significantly from the mean. Therefore, the approach is only applicable when there is a clear consensus between a reasonable number of reputation sources.

Dellarocas adopts a slightly different approach. First, he attempts to prevent negative discrimination through controlled anonymity, by which reputation and services are distributed by a central institution that does not reveal the identity of producers or consumers to each other. The intuition here is that, because a reputation source does not know the true identity of the trustee, it cannot determine if it is friend or foe and so has no reason to discredit it. This approach does not account for positive discrimination, because if a trustee and a reputation source collude, they could signal their identity to each other by other means, breaking anonymity.

To deal with this, the authors apply a clustering algorithm to separate a trustee's reputation into an upper and lower group of opinions. Since positive discrimination should appear more complimentary of the trustee, such opinions are assumed to be in the upper cluster, which is discarded. In most cases discarding the upper cluster introduces a negative bias. Through empirical study, the author argues that this bias is within acceptable bounds.

In our view, the main limitations of the Dellarocas approach lie in the applicability and effectiveness of controlled anonymity. Obviously, there are many cases in which a provider and consumer must be aware of each other's identity for a transaction to take place, which limits the situations in which this can be applied. Where it can be applied, it cannot account for a reputation source that wishes to discredit all trustees other than itself, or assumes that any agent that does not signal its true identity is a foe.

More generally, however, the assumption that majority opinion is reliable does not hold when there is a trustee with whom no agent has significant experience. In this case, all benevolent reputation sources will report no information, while reputation sources with an incentive to lie, will report information. In light of this, most, if not all, of the reputation provided will be unreliable.

⁷The interquantile range of a dataset is a descriptive statistic that specifies a range of values in which a given percentage of the data lie.

2.3.3 Past Performance

To alleviate these problems, we can consider the alternative exogenous approach of assessing a reputation source based on the accuracy of its past opinions. Amongst others, this is adopted by [Yu and Singh \(2003\)](#), who extend their previous work (Section 2.2.3) by applying a modified version of the Weighted Majority Algorithm ([Littlestone and Warmuth, 1994](#)). Essentially, their approach consists of three steps. First, the reputation of a trustee is calculated as a weighted average of reputation source opinions, with initially equal weights. Second, after the result of an interaction with the trustee has been observed, the differences between each opinion and the observed result are calculated. Third, the weights applied to each reputation source are adjusted relative to the difference between their stated opinion and the observed result.

There are two main advantages of this approach: First, under the reasonable assumption that a reputation source's past and future accuracy are correlated, the relative weight placed in inaccurate reputation sources will gradually decrease towards zero. Second, unlike [Dellarocas and Whitby](#), this approach does not require the majority of reputation opinions to be accurate, and so does not suffer the consequences associated with that assumption.

However, as we described previously, the method of representing trust employed by [Yu and Singh](#) does not adequately describe the uncertainty that surrounds an agent's behaviour. A better attempt at this is made by some probabilistic trust models, of which a good example is given by [Despotovic and Aberer \(2005\)](#). In their work, [Despotovic and Aberer](#) consider three different scenarios, two of which represent trustee behaviour as a binary event (cooperate or defect), and one which represents a trustee's performance during an interaction as a real value. In the latter case, the real value assigned to an agent's performance is supposed to represent some measure of its quality, and is assumed to be drawn from a normal distribution, with known variance, but unknown mean.

In both these scenarios, maximum likelihood estimation ([DeGroot and Schervish, 2002](#)) is used to derive the probability that a trustee will behave in a certain way, based on third party opinions. In particular, inaccurate opinions are handled by modelling the likelihood of receiving a certain opinion, given the probability that a reputation source is lying or telling the truth. In turn, the probability of a report being accurate is learnt over time, by a truster comparing reports about its own behaviour to its own knowledge of how it actually behaved.

The strength in this approach is that it shows how a trustee can be assessed in a well founded way, based on possibly inaccurate opinions. However, the assumption that a truster can learn about the reliability of reports based on its own reputation is not widely applicable. For instance, a particular agent might only ever consume a certain resource, and never provide it, and therefore never obtain a reputation for it. In addition, in two

of the three scenarios tested, the truster can only judge the average reliability of all reputation sources, as opposed to the accuracy of one individual. In this case, a truster is limited in how much it can improve its assessments, because it cannot distinguish between reliable and unreliable sources.

2.4 Mechanism Design

The techniques described so far have all addressed trust by attempting its assessment based on available knowledge. An alternative approach, *mechanism design* (Dash et al., 2003), aims to design a system in such a way that it is in the best interest of the agents to behave favourably towards each other. An established research area in its own right, this is not always explicitly tied to issues of trust, but from a trust perspective, it reduces the uncertainty surrounding a trustee's willingness to behave well. However, uncertainty in trustee behaviour cannot be removed completely. Generally, to manipulate a trustee's interests, we must assume that it is rational, which may not be the case for a variety of reasons, not least that an agent may have contracted a virus. Also, affecting an agent's willingness does not affect the uncertainty surrounding its capabilities. In light of this, we view mechanism design as complementary to trust assessment, rather than a replacement for it. Here, we illustrate how it can be used to simplify trust assessment problems, by reviewing some of the methods that lie in the intersection between trust and mechanism design.

Many of the trust models considered above include recommendations that can be considered as mechanism design. For instance, in HISTOS and SPORAS, it is not possible for an agent to have a reputation value lower than that of a new unknown agent entering the system for the first time. If this were possible, and agents were able to change identity at no cost, then agents with low reputation would simply create a new identity to improve their standing. Unfortunately, this approach may lead agents never to trust new agents if measures are not taken to ensure otherwise.

An important problem not addressed above is the incentive that an agent has to act as a reputation source. Clearly, if agents share information about trustee behaviour, they can increase their combined expected utility. However, this is not sufficient to persuade individual agents not to freeload, taking advantage of any available information, while not sharing any of their own. Jurca and Faltings attempt to alleviate this problem by introducing side payments for reputation (Jurca and Faltings, 2003). This obviously provides an incentive for an agent to supply reputation information, but it does not distinguish between accurate and inaccurate reputation. To rectify this, Jurca and Faltings suggest the following conditions should be guaranteed:

- Agents that report truthfully the result of every interaction with another agent, should not lose utility.

- Agents that report reputation incorrectly should gradually lose utility.

To ensure these conditions, Jurca and Faltings suggest that agents should only be paid for their opinion if it matches the next opinion received about the same trustee from a different source. Unfortunately, this approach fails if most agents provide false information, if agents collude to provide matching false reports, or if agents hold multiple identities to outwit the truster.

A more robust solution is provided by [Dash et al. \(2004\)](#), who introduce the concept of *trust-based mechanism design*, which attempts to explicitly handle issues of trust through mechanism design. In their approach, suppliers are allocated to consumers by a central institution (henceforth referred to as the centre). To aid the centre in making a good allocation, the consumer informs the centre of its preferences with regard to the allocation and all the information it currently knows about potential suppliers. Furthermore, the consumer either receives or makes a payment to the centre based on the effect its information has on the overall utility of the system. Based on these two components, it can be shown that it is in the best interest of a consumer to provide its reputation information fully and accurately.

One notable exception, however, is the possibility of agents colluding under certain conditions. A key premise is that agents will truthfully reveal their utility functions for an allocation, because to do otherwise risks decreasing the agent's utility in the allocation. This does not preclude the agent from omitting preferences it may hold that do not affect its allocation. For instance, suppose an agent wishes to decrease or increase the chances of another agent receiving a good allocation, and that it may further this goal by reporting inaccurate reputation. Provided the effect of this inaccuracy does not affect its own allocation, then it may do so without penalty, and for this reason, the approach only works when all of an agent's preferences concern its own allocation.

2.5 Assessing Unknown Agents

So far, we have seen three main approaches by which a truster can assess a trustee's behaviour:

1. It can draw upon its own personal experience or knowledge of a trustee.
2. It can seek the experience of other agents through reputation.
3. It can assume that a trustee will behave in a certain way, by making it in the trustee's best interest to do so.

However, there are cases where each of these approaches may fail. For instance, there may be little or no direct or indirect information about some trustees, and some types

of behaviour may not be assured by manipulating incentives. These types of issues are particularly apparent when we consider the problem of *whitewashing* (Zacharia *et al.*, 1999), whereby agents avoid punishment or bad reputation by assuming a new identity.

Nevertheless, even in these situations, a truster may need to make reasonable decisions about how to interact with such agents. One suggestion, made by Zacharia *et al.*, is that previously unknown agents entering a system should always be assigned the lowest possible level of trust. The main reason given for this is to discourage whitewashing, because it makes it impossible for an agent to improve its standing by assuming a new identity. The problem with this approach is that it unfairly discriminates against new members, who as a consequence may never establish themselves as useful members of the system.

As an alternative, we could assess unknown agents based on the behaviour of other similar agents for which some information is available. For whitewashing in particular, this offers a more pragmatic approach, by adapting a truster's decisions to best suit the general behaviour perceived in the environment. This type of solution is offered by REGRET and by Sun *et al.* (2005). However, in both of these cases, it is unclear precisely how much weight should be placed in this type of information over knowledge that is specific to a trustee itself. Determining how much predictive value group behaviour has in a given environment is thus an open question.

2.6 Summary

In this chapter, we have addressed three key points. First, we considered methods for representing and assessing trust based on evidence directly available to the truster. Second, we reviewed mechanisms for taking account of third party opinions, bearing in mind the extra challenges this source of trust imposes. Finally, we described the complementary role of mechanism design with regard to trust assessment; that is, how it can simplify the trust assessment problem, by reducing the uncertainty in a trustee's behaviour *a priori*. In this section, we summarise the main points made throughout the chapter, and identify key challenges for future research into trust assessment.

In Section 2.1 we considered work that concentrates on the cognitive aspects of trust — the core beliefs that a truster must hold to rationally be in a state of trust with a trustee. The main contribution of this work is that it helps to better understand the nature of trust, and the factors that contribute to it. However, it is not always clear how these core beliefs can be elicited from a truster's environment.

In contrast to this, the REGRET system demonstrates how a wide range of evidence can be brought together and used to assess trust in a given context. These sources include previous interactions with a trustee, third party opinions, information about

other agents in the same group as the trustee, the relationship between the truster and the trustee, and general assumptions about trustee behaviour. REGRET thus gives a reasonable assessment of a trustee, both when there is a significant amount of information available, and when information is scarce. The main disadvantage of REGRET is that it is based on *ad hoc* formulae, which require many parameter settings with no obvious optimal values.

The two main alternatives to *ad hoc* formulae, as found in REGRET, include Dempster Shafer theory, and probability theory. An example application of Dempster Shafer theory to trust is given by Yu and Singh who show how Dempster Shafer theory can be used to assess trust, based on previously observed interactions with a trustee. Although their method is sound in general, the way in which they ground trust in observed interactions is somewhat arbitrary.

Of the probabilistic trust models, the majority represent trust as the probability of a binary event; that is, the probability that a trustee will cooperate or defect. These models generally provide a sound statistical basis for calculating trust based on available evidence, and offer an attractive alternative to *ad hoc* formulae for trust assessment. However, by modelling a trustee's possible actions simply as cooperation or defection, they ignore the effect that the quality of service provided by a trustee may have on a truster. In addition, they remain dependent on (direct or indirect) observations of the trustee's own behaviour, and do not consider other sources of information, such as those explored by REGRET.

The vast majority of these trust models rely on third party opinions. Using such opinions, however, imposes several additional concerns that do not arise from knowledge directly available to the truster, because a third party's own preferences and world view inflict a bias on its opinions. To deal with this, there are three main approaches:

1. We can assess the reliability of opinions based on factors not pertaining to the opinions themselves.
2. We can assume that the majority of opinions are correct, and discredit opinions that deviate significantly from mainstream opinion.
3. We can assess the reliability of reputation, based on the perceived accuracy of reports received in the past.

Of these, we believe the third approach shows the most promise and potential for wide applicability, as it does not rely on external factors, nor assume that the majority opinion is correct. However, existing models that take this approach do not offer a perfect solution, either because they are founded on representations of trust that do not adequately capture the uncertainty surrounding a trustee's behaviour, or because they

assume a truster can determine a reputation source's accuracy based on information that is unlikely to be adequate for this task.

When a truster has neither reputation nor direct experience to draw upon, it has two other techniques that it can draw upon: either it can attempt to manipulate a trustee's incentives to make it in its best interest to behave in a certain way, or it can judge the trustee based on the behaviour of other similar agents in the system.

Of these, the former is realised through mechanism design, but does not offer a complete solution because not all types of behaviour can be ensured by manipulating incentives. On the other hand, the latter can provide useful information when no experience specific to the trustee is available. However, existing models that enable this approach do not adequately suggest the relative weight that such information should have against whatever little knowledge is available and specific to the trustee.

Against this background, we draw three main conclusions that motivate the work presented in the proceeding chapters. First, of all the existing methods for assessing trust, those that employ probability theory offer the strongest solution. This is because probability theory is well suited for representing and reasoning about the uncertainty surrounding events (such as agent behaviour), and is based on a sound theoretical foundation.

For this reason, in Chapter 3 we present a general framework for reasoning about trust based on probability theory, which forms the basis of two trust modes, TRAVOS and TRAVOS-C, detailed in Chapters 4 and 5 respectively. As discussed in Section 1.5, we have developed these models independently because different approaches are appropriate, depending on how trustee behaviour is represented. As a consequence, TRAVOS is designed for cases in which a trustee can only act in one of two ways during an interaction, while TRAVOS-C deals with cases in which trustees have a continuous range of possible behaviours.

Second, of the existing approaches for dealing with inaccurate reputation, those that assess the reliability of a reputation source based on past performance are the most widely applicable for the reasons stated above. However, as all of the models that adopt this technique leave room for improvement, we have developed methods in both TRAVOS and TRAVOS-C, which assess a reputation source based on past performance, while addressing some of the limitations of existing approaches. Specifically, TRAVOS includes a heuristic filtering mechanism that is shown empirically to outperform the method presented by Whitby *et al* (Section 2.3.2) in a number of important cases. Building on this, TRAVOS-C includes a more refined filtering method, which has a more solid theoretical foundation based on probability theory, without the need for heuristics.

Finally, a promising approach for making decisions when there is little specific information available about a trustee is to assess a trustee based on the behaviour of other agents in the system. However, assessing how much impact such information should

have on a truster's decisions remains an open question. To address this, in Chapter 5 we describe a probabilistic approach to this issue, as part of TRAVOS-C. Unlike the existing methods discussed in Section 2.5, this approach chooses the most appropriate estimate of a trustee's behaviour based on the behaviour of its peers, taking into account the amount of correlation observed between agents' behaviour in the environment.

Chapter 3

A Probabilistic Framework for Modelling Trust & Reputation

According to the objectives we laid down in Chapter 1, we wish to develop models of trust that draw upon both a truster's own experiences of a trustee, and the experiences of third parties. In the case of third party experience, however, we have three extra concerns: (1) that a benevolent third party should be able to convey all relevant information about its experiences to a truster; (2) that this should be achieved with minimum communication overhead; and (3) that where inaccuracies in reputation can occur, the truster's inferences should be robust against them.

With this in mind, the purpose of the current chapter is to provide a general framework for building trust models that fulfill these objectives. This not only acts as a foundation for the trust models we describe in later chapters, but also provides a set of notation that we shall use throughout the remainder of the thesis.

That said, questions of accuracy and communication are intimately tied to the semantics of the concepts we are trying to communicate. That is, if we wish to convey a concept or test its correctness, this can only be done with reference to its precise meaning. As such, our framework must be tied to a particular representation of trust, and given the conclusions of the previous chapters, this is a probabilistic one.

This chapter consists of eight parts: Sections 3.1 and ?? outline some basic concepts from decision theory and statistics, which we apply within the framework; Sections 3.2 and 3.3 define the basic problem we wish to tackle, along with notation, and show how this can be addressed using a truster's personal experience of a trustee; Sections 3.4 and 3.5 give guidelines for communication reputation between agents, in a way that meets our objectives in Chapter 1; finally, Section 3.6 categorises sources of inaccuracies in reputation, which must be dealt with by our trust models, and Section 3.7 summarises.

3.1 Background

In many situations, an agent must decide how best to act so as to achieve its goals. For example, in the context of this thesis, we want to determine if a truster should put its faith in a trustee to perform a given task. Moreover, if there are a number of competing trustees, we wish to determine which (if any) of these a truster should choose to interact with. As we have seen in the previous chapter, answering these sorts of questions does not require a complete reinvention of the wheel. Instead, we can turn to decision theory (Russell and Norvig, 2003), which defines a well established set of rules for solving these types of problems.

The key idea in decision theory is quite simple. Essentially, an agent should always choose actions that move it toward situations, or states, that it prefers. For example, suppose that an agent wants to build a house, and to do so it must choose between two builders, Dodgy Bill and Reliable Rod. The agent knows that if Dodgy Bill is awarded the building contract, he will take the money and run, whereas if Rod gets the contract, he will build the house fit for purpose. Clearly, the agent would prefer the state of owning a reliable house, to a state in which it has no house and no money. Thus, the best action is to award the contract to Rod, since this moves the agent toward the preferred state of owning the house.

To capture this intuition, an agent's preferences about the world are encoded in a *utility* function, which maps the set of situations the agent may find itself in to a set of real numbers representing the utility, or value, of each state to that agent. This is done such that, if one state is preferred to another, then the utility value returned for the preferred state will be higher than that of the other. Similarly, if two states are equally preferred, then they should have equal utility. Thus, to move to states that it most prefers, an agent should always choose actions that maximise its utility. Put formally, if S is the set of all states the agent can reach through its own actions, and the utility of each $s \in S$ is given by the utility function $U(s)$, then the agent should choose to act so as to arrive in a state s' , such that:

$$\forall s \in S, U(s) \leq U(s') \quad (3.1)$$

Unfortunately, knowing how to act to arrive in s' is not always feasible, and problems of trust are no exception. Typically, if a truster has to choose between competing trustees (such as Dodgy Bill and Reliable Rod) it will not know for certain which trustee will act most in its favour. To deal with this uncertainty, decision theory draws upon probability. Here, rather than maximise utility directly, an agent should act to maximise *expected utility*, defined as:

$$E[U(s)] = \int_S U(s)p(s|a) ds \quad (3.2)$$

where $p(s|a)$ is the probability density of s , given that the agent chooses action a . If this is adhered to, then an agent is not guaranteed to arrive in the best state every time. However, if the same situation were to be encountered n times, then the agent's average utility would approach its maximum as n becomes arbitrarily large. Essentially, this means that, although success is not always guaranteed, over many problems, the agent will perform to the best of its ability.

The presence of probability in the definition of expected utility means that decision theory in the face of uncertainty enters the realm of statistics. That is, to determine the expected utility of an action, we must determine the probability distribution of the possible states an action will result in. How we determine this distribution depends on our point of view, and the nature of the problem at hand. In particular, there are two main schools of thought in statistics, known as Bayesian and Frequentist statistics (DeGroot and Schervish, 2002), which differ over how a problem should be approached.

Generally speaking, Frequentists only consider probability with reference to the relative frequency of the different possible outcomes of an experiment which (at least conceptually) can be repeated a large number of times. On the other hand, Bayesians view probability as a subjective measure of the uncertainty surrounding a state, and are willing to assign probability to any statement, even if it does not depend on a random process. Moreover, in calculating a probability, a Bayesian routinely accounts for the subjective prior beliefs of the statistician in the calculation, whereas a Frequentist refuses to rely on anything but objective observations.

Many pros and cons have been argued for both stances, but this is not a debate that we need enter in detail here. Instead, for our purposes it is sufficient to adopt an eclectic stance, and draw on techniques from either side of the debate depending on what is convenient and appropriate. As we shall discuss in future sections, the problem of choosing trustees is one in which a truster may have little more than its own subjective beliefs to rely upon, or it may need to draw upon different sources of information with varying degrees of reliability. These factors can be much more conveniently handled using Bayesian rather than Frequentist methods, since these provide mechanisms for reasoning about subjective beliefs and also the uncertainty surrounding estimates of probabilities themselves. Therefore, for the rest of this section we focus on Bayesian estimate techniques that will be referred to later in the text.

Returning to the issue of expected utility, it is unusual for the state probability distribution to be of a known indisputable form. Instead, under both Frequentist and Bayesian philosophies, it is normally estimated given the assumptions and evidence available to the statistician. Typically, we assume that the shape of the distribution takes on one of a restricted set of forms, which is fully determined by an unknown parameter vector, θ . For example, suppose we wanted to determine the probability of obtaining heads from tossing a biased coin, having observed a set $X = \{x_1, \dots, x_n\}$ of previous tosses of the

coin. Here, the probability distribution of obtaining heads or tails would be a binary (or Bernoulli) distribution, for which θ could be the probability of obtaining heads¹. The distribution is then estimated either by estimating θ and substituting into Equation 3.2, or under the Bayesian framework, by marginalising over all possible values of θ .

In the case of estimating θ , this can be achieved by calculating its expected value using Equation 3.3, in which $p(\theta|X)$ is the probability density function (p.d.f.) of θ given the evidence. This approach has the desirable property of minimising the mean squared estimation error of θ , but it does not account for the amount of uncertainty surrounding its value. An alternative, although sometimes less tractable approach, is to marginalise over the space of possible values of θ , denoted Θ (Equation 3.4). Essentially, this averages the expected utility for all plausible values of θ . Thus, if there is little uncertainty about θ , the result will be similar to that of the true θ . On the other hand, if there is much uncertainty about θ , then there will be a larger number of plausible values for θ having a significant impact. This has the effect of making the result less susceptible to error than accounting only for an estimated θ that may deviate significantly from the true value:

$$E[\theta] = \int_{\Theta} \theta p(\theta|X) d\theta \quad (3.3)$$

$$E[U(s)] = \int_S U(s) \int_{\Theta} p(s|a, \theta) p(\theta|X) d\theta ds \quad (3.4)$$

In either case, we are left with the need to determine the posterior density of θ given the evidence, which (assuming each $x_i \in X$ is independent) is achieved by applying Bayes rule as shown in Equation 3.6. Here, $p(X|\theta)$ is known as the data likelihood function, which usually follows directly from the model assumptions. For example, in the case of our biased coin tosses, the likelihood of each x_i is $p(x_i|\theta) = \theta$ if x_i is heads, or $p(x_i|\theta) = (1 - \theta)$ otherwise.

With regard to the other factors, $p(\theta)$ is the prior density of θ , which is chosen in line with the prior beliefs or assumptions of the statistician. For example, if prior to observing any coin tosses all values of θ are believed equally likely, then $p(\theta)$ can be chosen to be a uniform distribution. Such a prior is appropriate if the statistician has no information about θ *a priori*, and often has little impact on the posterior distribution given even a small amount of data.

On the other hand, existing information from other sources may be incorporated into the prior, to bias estimates towards values of θ believed to be more likely. For example, if θ is a parameter relating to the quality of service provided by a grid computing resource, then a prior distribution could be used that is based on some prior modelling of factors pertaining to the application domain. However, such modelling is beyond the scope of

¹This fully specifies the distribution since the probability of obtain heads is then θ , while the probability of obtaining tails is $(1 - \theta)$.

our work.

$$p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} \quad (3.5)$$

$$= \frac{p(\theta) \prod_{i=1}^n p(x_i|\theta)}{p(X)} \quad (\text{assuming all } x_i \text{ are independent}) \quad (3.6)$$

Finally, $p(X)$ is the marginal distribution of the data. Essentially, this can be viewed as a normalising constant that ensures that $\int_{\Theta} p(\theta|X) dX = 1$, as should be the case for any p.d.f. As such, the value of $p(X)$ is determined by integrating the numerator in Equation 3.6 over the domain of θ , which along with the integrals in Equations 3.3 and 3.4 often does not result in a closed form analytical solution. This leaves us with two choices: either we constrain our model to cases that do have analytical solutions, or we turn to numerical integration techniques (Appendix A). Constraining the model may not always be appropriate, but in some cases, this can be achieved by choosing a particular form for the prior, which in many cases may not be too restrictive. To give a concrete example, suppose that we have the following prior p.d.f. for the bias of our coin:

$$p(\theta) = \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{Beta(\alpha, \beta)} \quad \text{where,} \quad (3.7)$$

$$Beta(\alpha, \beta) = \int_0^1 t^{\alpha-1}(1-t)^{\beta-1} dt \quad (3.8)$$

Distributions with this density are known as beta distributions, in which α and β are shape parameters, and $Beta(\alpha, \beta)$ is the beta function, which is defined in Equation 3.8. Now, if we toss the coin a number of times and observe a total of N heads and M tails then, from Bayes rule, the posterior distribution is proportional to:

$$p(\theta|X) \propto p(\theta) \prod_{i=1}^l p(x_i|\theta) \quad (3.9)$$

$$\propto \theta^{\alpha+N-1}(1-\theta)^{\beta+M-1} \quad (3.10)$$

$$\therefore p(\theta|X) = \frac{\theta^{\alpha+N-1}(1-\theta)^{\beta+M-1}}{Beta(\alpha+N, \beta+M)} \quad (3.11)$$

Thus, by substituting the original shape parameters for $\alpha + N$ and $\beta + M$ respectively, we see that, like the prior, the posterior distribution is also a beta distribution. Prior distributions that exhibit this property of ensuring the posterior shares its form are known as conjugate priors (see Section 2.2.2). They are particularly significant if we know existing properties, such as how to calculate the density, because this means the same results apply to the posterior as apply to the prior. In this case of the beta distribution, we do know how to efficiently calculate its p.d.f. and expected value, which greatly simplifies the estimation problem.

Another example that we will refer to in later chapters is Bayesian inference concerning Gaussian (or normal) distributions. For this, if we have l i.i.d. samples $X = \{x_1, \dots, x_l\}$, which are drawn from a Gaussian with unknown mean (μ) and variance (σ^2), then from the Gaussian p.d.f. we have the following likelihood functions:

$$p(x_i|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right] \quad (3.12)$$

$$p(X|\mu, \sigma) = \prod_{i=1}^l \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right] \quad (3.13)$$

For this, the conjugate prior can be expressed in terms of $p(\mu|\sigma)$ and $p(\sigma)$ as:

$$p(\mu, \sigma^2) = p(\mu|\sigma^2)p(\sigma^2) \quad (3.14)$$

$$p(\mu|\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2/\lambda}} \exp\left[-\frac{(\mu - m)^2}{2\sigma^2/\lambda}\right] \quad (3.15)$$

$$p(\sigma^2|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} \exp\left[-\frac{\beta}{\sigma^2}\right] \quad (3.16)$$

where $p(\mu|\sigma)$ is a Gaussian density with mean m and variance σ^2/λ , and $p(\sigma)$ is an inverse-gamma distribution with shape parameters α and β . From this, it can be shown that the posterior parameter distribution has the same form, but with corresponding hyperparameters, μ' , λ' , α' and β' , derived as follows (see [DeGroot and Schervish \(2002\)](#) for details).

$$\mu' = \frac{\lambda\mu + l\bar{x}}{\lambda + l} \quad (3.17)$$

$$\lambda' = \lambda + l \quad (3.18)$$

$$\alpha' = \alpha + \frac{l}{2} \quad (3.19)$$

$$\beta' = \beta + \frac{s^2}{2} + \frac{l\lambda(\bar{x} - \mu)^2}{2(\lambda + l)} \quad \text{where,} \quad (3.20)$$

$$\bar{x} = \frac{1}{l} \sum_{i=1}^l x_i \quad (3.21)$$

$$s^2 = \sum_{i=1}^l (x_i - \bar{x})^2 \quad (3.22)$$

3.2 Basic Notation and Problem Definition

So far, we have described how decision theory can be used in conjunction with probability theory to enable an agent to choose between possible actions. We now wish to show how this general theory can be applied to our specific problem, where trusters need to choose whether or not to interact with a potential trustee. However, to see how we

approach this problem, we must first introduce some basic notation, which we shall reuse throughout the rest of the thesis.

To this end, in a MAS consisting of n agents, we denote the set of all agents as $\{a_1, a_2, \dots, a_n\} = \mathcal{A}$. Over time, interactions take place between distinct pairs of agents from \mathcal{A} , during which one of these agents is obliged to provide a service to the other. In each case, the agent receiving the service is the truster, denoted a_{tr} , and the agent providing the service is the trustee, denoted a_{te} .

With an aim to assess trustee performance, a truster records the outcome of each interaction as it *perceives* it, which is denoted as $O_{a_{tr}, a_{te}}$ — the outcome of interacting with a_{te} from the perspective of a_{tr} . From this interpretation, bilateral interactions in which both parties have obligations to each other can be seen as two separate interactions in which each agent plays the role of truster and trustee in turn. If such an event occurs between agents a_1 and a_2 , then this will result in two recorded outcomes, denoted O_{a_1, a_2} and O_{a_2, a_1} . However, it is important to note that O_{a_1, a_2} and O_{a_2, a_1} are not necessarily equal, as each agent may represent the outcome only in terms that are relevant to it. For example, if a_1 sells high quality apples to a_2 , for which a_2 does not pay, then from a_2 's perspective the interaction results in the possession of some high quality apples, while from a_1 's perspective, goods are lost without payment.

With this in mind, it is useful to define a number of outcome instances, and sets involving them. First, we define the set of all possible outcomes in a particular context, \mathcal{C} , as $\mathcal{O}^{\mathcal{C}}$. Here, a context specifies both the type of interaction from which outcomes are derived and the way it is recorded. For instance, in the example given above, we could have $O_{a_2, a_1} \in \mathcal{O}^{apples}$ and $O_{a_1, a_2} \in \mathcal{O}^{money}$, where each context is defined in terms of the services received by the respective truster.

In general, the rest of our discussion in this and following chapters applies independent of context. Nevertheless, when assessing a trustee's behaviour in a given context, we only base our predictions based on outcomes in that context, as a trustee's behaviour in one context does not necessarily provide information about its behaviour in another.² Therefore, when we do discuss sets of outcomes and functions defined on them, we assume that all such outcomes belong to the same context, regardless of what that context actually is. The context superscript for outcome spaces acts as a reminder of this.

Building on this, we divide time into discrete steps starting from time 0, and denote the outcome of an interaction that occurred between a_{tr} and a_{te} at time t as $O_{a_{tr}, a_{te}}^t$. In general we wish to allow any number of interactions to occur between any agents at any time. However, to simplify our discussion we will assume that at most one interaction can

²Although a trustee's behaviour in one context may provide information about its behaviour in another, using such information would require modelling the dependences that exist between contexts, which is beyond the scope of our work.

occur between a given truster and trustee in a given time step, and that each interaction is complete by the end of the time step it is said to occur in. Furthermore, we denote the current time as t' , and the set of all outcomes between a_{tr} and a_{te} from time t to $t+n$ as $O_{a_{tr},a_{te}}^{t:t+n}$. Thus, the history of all interactions between a_{tr} and a_{te} is given by $O_{a_{tr},a_{te}}^{0:t'}$.

3.3 Trust Assessment Based on Direct Observations

Now that we have a formal language for discussing interactions between agents, the key question is how can we apply decision theory to the assessment of potential trustees in this context. The answer to this comes in two parts. First, we assume that the utility received by a truster for interacting with a trustee is completely determined by the outcome of that interaction. This means that the preferences of any a_{tr} with regard to interacting with a_{te} can be encoded by a utility function $U : \mathcal{O}^C \rightarrow \mathbb{R}$, such that if $O_{a_{tr},a_{te}}^{(1)}$ is preferred over $O_{a_{tr},a_{te}}^{(2)}$, then $U(O_{a_{tr},a_{te}}^{(1)}) > U(O_{a_{tr},a_{te}}^{(2)})$; and if $O_{a_{tr},a_{te}}^{(1)}$ is equally preferred to $O_{a_{tr},a_{te}}^{(2)}$, then $U(O_{a_{tr},a_{te}}^{(1)}) = U(O_{a_{tr},a_{te}}^{(2)})$. Second, we assess the value of a_{tr} interacting with a_{te} by calculating the expected utility according to Equation 3.23 (assuming \mathcal{O}^C is continuous) or Equation 3.24 (assuming \mathcal{O}^C is discrete).

$$EU = \int_{\mathcal{O}^C} U(O_{a_{tr},a_{te}}) p(O_{a_{tr},a_{te}}) dO_{a_{tr},a_{te}} \quad (3.23)$$

$$EU = \sum_{O_{a_{tr},a_{te}} \in \mathcal{O}^C} U(O_{a_{tr},a_{te}}) p(O_{a_{tr},a_{te}}) \quad (3.24)$$

The precise definition of $U(O_{a_{tr},a_{te}})$ in this equation is something that depends on the particular application at hand, and so is not something that we address here. Calculating $p(O_{a_{tr},a_{te}})$ is, however, something that can be discussed in more general terms, and is what we wish to address through our models of trust. How this is achieved depends on the evidence used, but here we shall concentrate on the history of trustee observations $O_{a_{tr},a_{te}}^{0:t'}$, and a_{te} 's reputation.

Of these, trust based on direct experience is the most straightforward, following directly from standard practice as described in Section ???. Specifically, we assume that all outcomes of interactions between a_{tr} and a_{te} are independently drawn from a single distribution, which is fully determined by a parameter vector $\theta_{a_{tr},a_{te}}$ with domain Θ^C . In essence, this distribution summarises the uncertainty that is intrinsic to the behaviour of a_{te} toward a_{tr} , and as such will be referred to as the trustee's *behaviour* distribution. For example, if this distribution has relatively high variance, then this means that the behaviour of the trustee is relatively unpredictable. On the other hand, if the distribution is highly peaked around a single value with low variance, then this means that the behaviour of the trustee is consistent and highly predictable.

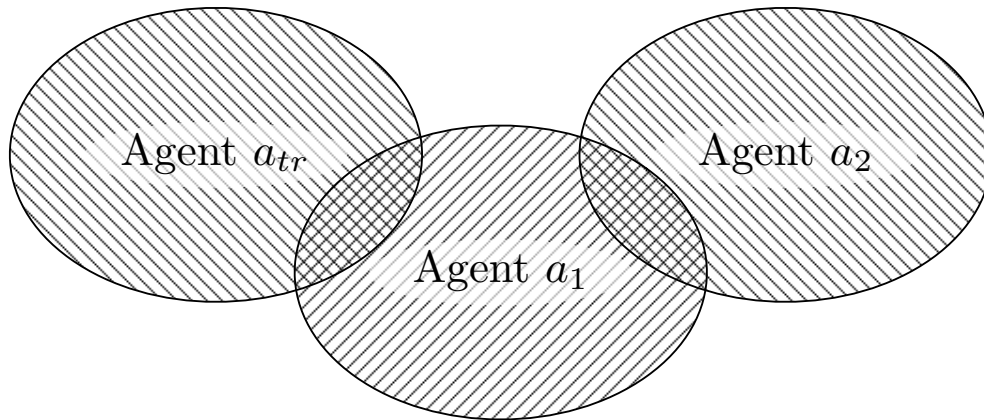


FIGURE 3.1: Venn diagram of overlapping reputation datasets.

Nevertheless, this is not the whole story because there is likely to be an amount of uncertainty surrounding the value of $\theta_{a_{tr},a_{te}}$ depending on the amount of evidence or knowledge available. This we characterise in the manner described previously, by modelling the parameter distribution of $\theta_{a_{tr},a_{te}}$. Accounting for a truster's experience of a trustee is then just a matter of obtaining the posterior distribution of $\theta_{a_{tr},a_{te}}$ by applying Bayes rule as follows:

$$p(\theta_{a_{tr},a_{te}} | O_{a_{tr},a_{te}}^{0:t'}) = \frac{p(O_{a_{tr},a_{te}}^{0:t'} | \theta_{a_{tr},a_{te}}) p(\theta_{a_{tr},a_{te}})}{p(O_{a_{tr},a_{te}}^{0:t'})} \quad (3.25)$$

$$= \frac{p(\theta_{a_{tr},a_{te}}) \sum_{o \in O_{a_{tr},a_{te}}^{0:t'}} p(o | \theta_{a_{tr},a_{te}})}{p(O_{a_{tr},a_{te}}^{0:t'})} \quad (3.26)$$

In calculating this posterior, we may use any one of a number of techniques, but where we use conjugate priors, some extra notation applies. Specifically, the parameter distribution belonging to a trustee's behaviour is characterised by a hyperparameter vector $\phi_{a_{tr},a_{te}}$ with domain Φ^C . In this case, the posterior update is realised by the particular update rule that is associated with the parameter model being applied, and is defined in general by:

$$p(\theta_{a_{tr},a_{te}} | \phi_{a_{tr},a_{te}}^{post}) = \frac{p(\theta_{a_{tr},a_{te}} | \phi_{a_{tr},a_{te}}^{prior}) \sum_{o \in O_{a_{tr},a_{te}}^{0:t'}} p(o | \theta_{a_{tr},a_{te}})}{p(O_{a_{tr},a_{te}}^{0:t'} | \phi_{a_{tr},a_{te}}^{prior})} \quad (3.27)$$

where $\phi_{a_{tr},a_{te}}^{prior}$ and $\phi_{a_{tr},a_{te}}^{post}$ are the prior and posterior hyperparameters respectively.

3.4 Reputation Framework

The next problem we wish to tackle is how to perform trust assessment based on both a truster's direct experiences, and the trustee's reputation. In particular, we are interested

in reputation based on third party observations of a trustee's behaviour that, in the simplest of cases, can be achieved by generalising the technique applied in the previous section. For instance, suppose that a_{tr} has three (possibly overlapping) datasets on which to base its assessment of a_{te} : its own dataset $O_{a_{tr},a_{te}}^{0:t'}$, and the dataset of two other agents, a_1 and a_2 (Figure 3.1). If the truster had all relevant details about these datasets available, then it could account for all the evidence simply by substituting the union $O_{a_{tr},a_{te}}^{0:t'} \cup_{i=1}^2 O_{a_i,a_{te}}^{0:t'}$ for $O_{a_{tr},a_{te}}^{0:t'}$ in Equation 3.26.

However, in many cases we do not expect this to be a viable solution for two reasons. First, there may be an overhead associated with communicating reputation between agents. Hence, transmitting an agent's observations in their entirety may prove to be a costly enterprise. Second, the observations in $O_{a_i,a_{te}}^{0:t'}$ may not be drawn from the intended behaviour distribution, at least as they are reported. This may be because the trustee's behaviour differs depending on its interaction partner, resulting in $\theta_{a_{tr},a_{te}} \neq \theta_{a_i,a_{te}}$, or because a_i intentionally manipulates its reports to mislead the truster.

Although both of these points are associated with the objectives we outlined in Chapter 1, we begin by focusing our attention on the first. This we frame by assuming each reputation source, $a_{rep} \in \mathcal{A}$, has a function r , such that $R_{a_{rep},a_{te}} = r(O_{a_{rep},a_{te}}^{0:t'})$, where $R_{a_{rep},a_{te}}$ is the *opinion* of a_{rep} about a_{te} , and r is the *opinion function*. Also, in the interest of simplicity, we assume there is one shared definition of r for all agents, and that the datasets on which each agent bases its reported opinion do not intersect.³

Intuition then tells us that, rather than encode an agent's entire observation set verbatim, the opinion function should somehow capture all relevant information about $O_{a_{rep},a_{te}}^{0:t'}$, yet as concisely as possible. Of course, what information is relevant, depends on what it is that the information is supposed to be relevant to. In our case, we are trying to locate the true value of $\theta_{a_{tr},a_{te}}$, so are only interested in properties of an agent's experiences that tell us something about that parameter.

With this in mind, we turn to the concept of sufficient statistics, which captures this intuition by definition. For example, suppose that we have a set of i.i.d. variables $X = \{x_1, \dots, x_n\}$ that are drawn from a distribution with p.d.f. $f(x|\theta)$, where θ is a parameter which we wish to estimate. Any real-valued function $r = s(X)$ defined on such a set of observations is called a *statistic* (Definition 3.1). Moreover, if $r = s(X)$ is a sufficient statistic, with respect to θ , then the distribution of X given r does not depend on θ (Definition 3.2).

Sufficient statistics have many applications regarding parameter estimation, and can be found in a number of ways, most notably by applying the *factorisation criterion* as defined in Theorem 3.3. Our interest, however, stems from the relationship between

³It may be possible to loosen these assumptions if agents communicate (or estimate) intersections between their observations, and differences between their opinion functions. However, dealing with such cases is beyond the scope of our work.

sufficient statistics and the distribution of a parameter. Specifically, it follows from Theorem 3.4 that the distribution of θ given r does not depend on X .

Thus, if we have a sufficient statistic for a parameter, then no extra knowledge about a sample will ever affect the parameter's distribution, and so it will never affect our inferences about that parameter. Conversely, if we don't have a sufficient statistic for a parameter, then there will be some extra information about a sample that will improve our estimates about a parameter, if that information were made available.

Definition 3.1 (Statistic). Assume that X is a set of random variables, with domain \mathcal{X} , that corresponds to a set of observations. Then, a statistic of X is any function $s(X)$ that is defined for \mathcal{X} (adapted from Upton and Cook (2002)).

Definition 3.2 (Sufficient Statistic). If $r = s(X)$ is a statistic of a sample X , then $s(X)$ is said to be sufficient for a parameter θ , if and only if the probability distribution of X , given r , does not depend on θ . Equivalently, $T(X)$ is a sufficient statistic for θ if Theorem 3.3 holds (adapted from DeGroot and Schervish (2002)).

Theorem 3.3 (The factorisation criterion). Let X_1, \dots, X_n form a random sample from either a continuous distribution or a discrete distribution for which the probability density function (p.d.f.) or the probability function⁴ (p.f.) is $f(x|\theta)$, and where the value of θ is unknown and belongs to a given parameter space Θ . A statistic $T = r(X_1, \dots, X_n)$ is a sufficient statistic for θ if and only if the joint p.d.f. or joint p.f., $f_n(\mathbf{x}|\theta)$ of X_1, \dots, X_n , can be factored as follows for all values of $\mathbf{x} = (x_1, \dots, x_n) \in R^n$ and all values of $\theta \in \Theta$:

$$f_n(\mathbf{x}|\theta) = u(\mathbf{x})v[r(\mathbf{x}), \theta]. \quad (3.28)$$

Here, the functions u and v are nonnegative; the function u may depend on \mathbf{x} but does not depend on θ ; and the function v depends on θ , but depends on the observed value of \mathbf{x} only through the value of the statistic $r(\mathbf{x})$.

Theorem 3.4 (Parameter Distributions from Sufficient Statistics). Suppose that $X = \{x_1, \dots, x_l\}$ is a set of l samples independently drawn from a distribution with parameter θ . Then, if $r = s(X)$ is a sufficient statistic of X , the distribution of θ will be conditionally independent of X , given r ; put another way, the following equality is true.

$$p(\theta|X, r) = p(\theta|r) \quad (3.29)$$

⁴A probability function fulfills the same role for discrete probability distributions, as a probability density function does for continuous distributions (DeGroot and Schervish, 2002).

Proof: Using Bayes rule, along with Definition 3.2, we can derive the conditional density of θ given X and r as follows.

$$p(\theta|X, r) = \frac{p(X, r, \theta)}{p(X, r)} \quad (\text{Bayes rule}) \quad (3.30)$$

$$= \frac{p(X|r, \theta)p(r, \theta)}{p(X|r)p(r)} \quad (3.31)$$

$$= \frac{p(X|r)p(r, \theta)}{p(X|r)p(r)} \quad (\text{from Definition 3.2}) \quad (3.32)$$

$$= \frac{p(r, \theta)}{p(r)} \quad (3.33)$$

$$= p(\theta|r) \quad (\text{Bayes rule}) \quad (3.34)$$

Definition 3.5 (Minimal Sufficient Statistic). A statistic $r = s(X)$ is a *minimal sufficient statistic* of a sample X if r is a sufficient statistic X and is a function of every other sufficient statistic of X (DeGroot and Schervish, 2002).

In this light, we can conclude that the opinion function, r , should ideally be defined as a sufficient statistic, but this, on its own, does not guarantee that the opinion will be as concise as possible. We thus further strive to make the opinion function a *minimal* sufficient statistic (Definition 3.5), which ensures that only information that affects the parameter distribution is retained. This is not a hard constraint, as any sufficient statistic will still retain all relevant information about a sample. However, if a sufficient statistic is minimal, it will express the relevant information using the smallest space of possible values, making it easier to communicate the information efficiently.

To see why this works, suppose that Alice tosses a coin and wants to inform Bob of the result. To do so, she can choose between two possible codes: either she can transmit 1 for heads or 0 for tails, or she can transmit a number between 1 and 8, with the numbers 1 to 4 indicating heads, and 5 to 8 indicating tails. The difference between these codes is that there is a one-to-one relationship between the toss outcome and the binary code, but a one-to-many relationship between the outcome and the alternative code. Both codes distinguish between the states we wish to communicate, but the one-to-many relationship includes many redundant states that take more bits to represent and transmit.

Similarly, when we transmit a statistic about a data set, we are interested in what the posterior parameter distribution should be. If the statistic is sufficient, then the number of possible parameter distributions that may result is as large as can be. However, only if the statistic is minimally sufficient, will it have a one-to-one relationship with the space of parameter distributions, and so have the potential to enjoy the most concise representation possible, without loss of information.

3.5 A Word of Warning on Sufficiency

Although the discussion in the previous section realises all of our objectives with regard to communicating reputation, there are three important points about these recommendations that must be kept in mind.

First, we emphasize that these are only guidelines for getting the most value out of reputation as possible. Even if an opinion function is not sufficient, this does not necessarily mean that it provides no useful information. In addition, choosing a minimally sufficient statistic does not guarantee concise enough communication in every case, so it may be appropriate to choose a non-sufficient statistic to further reduce communication overhead.

Instead, minimal sufficiency should be seen as an optimal point in the trade-off between communication overhead and maintaining information: on one hand, anything that has a larger representation than a minimally sufficient statistic is wasting space somewhere, whether it is sufficient or not, while on the other hand, if we require a representation smaller than that of a minimal sufficient statistic, then we should be aware of how much information we are losing in the process.

Second, the definition of sufficiency is intimately tied up with the particular parameter model being used, so sufficiency for one model does not imply sufficiency for another. For example, suppose that we draw samples from a Gaussian distribution that we assume has a variance of 1. In this case, provided we know the size of the sample, the sample mean is sufficient for the model. However, once we remove this assumption, the extra free parameter means that a statistic requires more information to be sufficient than the sample mean can provide.

Finally, we do not expect the number of observations that a reputation source has of a trustee to be directly observable by the truster. This means that, in order to be sufficient, the sample size must be included, or be derivable from, the opinion function. We highlight this because there are instances when it is implicitly assumed in the literature that the sample size is observable.

For instance, in [DeGroot and Schervish \(2002\)](#) the sample mean is quoted in an example as being sufficient for estimating the mean of a Gaussian distribution. However, we know from the discussion in Section ?? that, when we use conjugate priors, the posterior distribution of a Gaussian mean, given the variance, does depend on the size of the sample. In addition, the same can be said for any model for which larger samples provide more information about the model parameters. As this is generally the case, the number of observations should be considered an essential part of most opinion functions.

3.6 Coping with Inaccurate Reputation

From the previous section, we have a fully specified framework for assessing trust based on reputation. However, for this framework to provide reasonable results for a given truster-trustee pair, then the following conditions must hold:

Condition 1. If $a_{tr} \in \mathcal{A}$ is a truster and $R \subseteq \mathcal{A}$ is the set of all reputation sources that a_{tr} consults about a trustee $a_{te} \in \mathcal{A}$, then the behaviour of a_{te} towards all members of $\{a_{tr}\} \cup R$ must be equal.

Condition 2. If $a_{tr} \in \mathcal{A}$ is a truster and $R \subseteq \mathcal{A}$ is the set of all reputation sources that a_{tr} consults about a trustee $a_{te} \in \mathcal{A}$, then all members of R must report their information about a_{te} truthfully and accurately.

Essentially, Conditions 1 and 2 ensure that observations made by a truster's reputation sources are representative of the actual behaviour a trustee is likely to have towards the truster. Unfortunately, we cannot expect these conditions to hold in general, so we must develop methods for coping with cases in which they are violated. Many of the trust models we reviewed in Chapter 2 include methods for coping with some of these conditions. However, as we stated in Section 2.6, each has its own set of problems.

To address these limitations, each of the trust models described in the following chapters include their own unique mechanisms for coping with the problems that arise when Conditions 1 or 2 are violated. Although we defer discussing possible solutions to these problems until then, it is useful to be aware of the types of issues that any such solution must account for.

To this end, the subsections that follow categorise the possible causes of inaccuracies in reputation, and define the impact that each category may have, with respect to the communication guidelines set forth in the previous section. In particular, we identify three groups of causes:

1. statistical noise, due to a trustee's own behaviour;
2. the opinion view of an agent, which defines how it perceives a trustee's behaviour;
and
3. the incentives of an agent, which may lead it to manipulate a trustee's reputation.

3.6.1 Statistical Noise

Statistical noise is noise which we normally associate with learning about a distribution by sampling from it under an i.i.d. assumption, and it is entirely due to the properties of the behaviour distribution. This type of noise does not affect the bias, but does affect

the error variance in that the error variance will decrease as the amount of information the sample conveys about the distribution increases. In particular, under reasonable assumptions, the error variance decreases as the sample size increases.

3.6.2 Opinion View

We define a reputation source's opinion view as its view of a trustee's behaviour; that is, how a trustee behaves towards the reputation source, and how that reputation source observes and records the trustee's behaviour. The opinion view of a reputation source, a_{rep} , is a source of noise in observations sent to a a_{tr} if it causes those observations to be drawn from a distribution other than $p(O_{a_{tr},a_{te}}|\theta_{a_{tr},a_{te}})$. For example, this can be the case if observing a trustee's behaviour involves taking readings from physical sensors, and the sensors used by a_{rep} are uncalibrated with those used by a_{tr} . On the other hand, the behaviour of the trustee may differ depending on whether it is interacting with a_{tr} or a_{rep} . This could be for environmental reasons, for example differences in network bandwidth used by a_{tr} and a_{rep} , or simply because the trustee has incentives to behave differently towards different agents.

As opinion views can cause observations to be drawn from different unknown distributions, they can affect both the bias and error variance of a truster's estimate. However, this is not to say that if the distributions differ, a third party's observations carry no information about $\theta_{a_{tr},a_{te}}$ because the distributions may still be correlated. For example, suppose that a_{te} provides video streaming to both a_{rep} and a_{tr} , but the networks connecting a_{te} to a_{rep} and a_{te} to a_{tr} offer different levels of service. Although in this case the network may cause different behaviour to be observed by a_{rep} and a_{te} , if the servers hosting a_{te} 's video content continually go off-line, this effect will be observed by both parties.

3.6.3 Opinion Incentives

When a truster assesses a trustee based on observations reported from third parties, the truster must consider the possibility that the reputation source may not reveal those observations truthfully. To deal with this possibility, the truster should consider any knowledge it has about the incentives the reputation source may have for, or against, truthtelling, and the ways in which the observations may be manipulated to meet such goals. Although a reputation source may have conflicting incentives, the net effect will fall into one of three categories:

Truthtelling Incentive In this case, the truster knows that it is in the best interest of the opinion provider to reveal information about its past experiences with the

trustee truthfully. Thus, the incentive of the opinion provider does not add any extra noise over and above that described in the previous sections.

Competitive Incentive Here, the truster believes that an opinion provider may wish to bias the truster’s decisions to the disadvantage of the trustee. Assuming that the opinion provider knows which types of trustee behaviour the truster dislikes, the following effects may occur:

- The opinion provider’s reported observations may be biased towards predicting unfavourable trustee behaviour, which we call a *negative* bias.
- If the opinion source is lying, the sample variance of the observations may be affected depending on the adopted lying strategy, because the samples may no longer be drawn from the actual behaviour distribution. Moreover, the lying strategy will also determine how bad the effect is. For example, the reputation source could completely discard its knowledge of the trustee’s true behaviour, in which case its observations may be completely uncorrelated with $p(O_{a_{tr},a_{te}}|\theta_{a_{tr},a_{te}})$. On the other hand, it could mediate its lies based on its knowledge to purposely introduce some correlation to disguise its lies. In the latter case, the correlation still provides useful information (albeit possibly reduced), even though it is introduced only to diminish the footprint of the reputation source’s lying behaviour.
- The opinion provider may exaggerate the sample size by reporting an observation set which has higher cardinality than any true underlying sample. This is because increasing the proportion of O_u that the provider is responsible for increases the provider’s influence over the truster’s final estimate. On the other hand, the opinion provider is unlikely to report the cardinality as being lower than its true value, because it is unlikely to have an incentive to decrease its influence on the truster’s opinion. In any event, even if the opinion provider does report a lower cardinality, the truster would be unwise to assume the cardinality to be higher than reported, because it would then be making estimates based on knowledge that it does not have.

This last effect is perhaps the most damaging consequence of reputation source incentives conflicting with those of a truster. With the noise sources discussed in Sections 3.6.1 and 3.6.2, we can at least be sure that the number of observations is correct, since there is known uncertainty surrounding this from the perspective of the observer. With conflicting interests, the truster must also consider the possibility that a set of reported observations is not even a real random sample, but a set concocted from the mind of the reputation source.

Collusion Incentive Here, the truster believes that an opinion provider may wish to bias the truster’s decisions to the *advantage* of the trustee. The effects of this source of inaccuracy, and the reasoning behind it, are the same as those of competitive

incentives, but in this case we assume that the effect on the bias is *positive* rather than negative. That is, assuming that the opinion provider knows which types of trustee behaviour the truster dislikes, the opinion provider may manufacture its reported observations so as to bias the truster's decisions in favour of the trustee.

3.7 Summary

In this chapter, we have considered problems in which an agent must choose whether or not to interact with a potential trustee, and have made three main points about how such problems may be solved.

First, solving such problems does not require a completely new solution, but instead we can draw upon existing techniques in decision theory and statistics. That is, when deciding whether or not to interact with a trustee, a truster should compare the interaction to its alternatives, and choose the course of action that maximises its chance of obtaining an outcome that it prefers.

To achieve this, the truster must have defined a utility function, which encodes its preferences about the possible outcomes of each course of action. The truster can then account for the uncertainty surrounding such outcomes, by using this function, along with the probability distribution of possible outcomes, to calculate the expected utility of each course of action, and then act in way that maximises this value. Moreover, to derive the probability distribution of possible outcomes, the truster can apply Bayesian analysis to its past experiences of a trustee (if any), along with reported third party experiences, known as reputation.

Second, if we wish to make decisions that account for a trustee's reputation, we must decide how best to communicate third party experiences between agents. Ideally, this should be achieved by communicating all relevant information as concisely as possible, which can be accomplished by defining an agent's opinion about a trustee as a minimal sufficient statistic of the agent's experiences with that trustee. This ensures that, aside from its reported opinion, a reputation source cannot convey any other information about its experiences that could improve a truster's assessment of a trustee. Conversely, by ensuring the statistic is minimal, we ensure that no effort is wasted transmitting information that does not have predictive value with regard to a trustee's behaviour.

Finally, in all but the most restrictive conditions, third party experiences cannot be assumed to be as reliable as a truster's own direct observations for two main reasons:

1. A trustee's behaviour, as it is measured, may be different depending on which agent it interacts with.

2. A truster's reputation sources may have incentives to misrepresent their experiences, in pursuit of their own goals.

Together these imply that, while reputation may generally provide useful information about a trustee, it may be noisy, or contain no useful information about a trustee at all.

In the following chapters, we apply the principles described in this chapter to cases in which trustee behaviour is represented in one of two ways. Specifically, Chapter 4 shows how, within TRAVOS, Bayesian analysis is applied to cases in which trustee behaviour is represented as a binary event, while Chapter 5 shows how the same is achieved by TRAVOS-C, when trustee behaviour is given a continuous representation. What is more significant, however, is the way in which each of these models deals with reputation. In this regard, we shall discuss how each model applies the recommendations of the current chapter, by ensuring efficient communication of reputation, and by addressing the possibility of inaccurate reputation.

Chapter 4

TRAVOS: A Trust Model for Boolean Action Spaces

So far we have seen that, when a truster and trustee interact, the worth of the interaction to the truster is determined by the actions of the trustee. This means that to determine what a truster stands to gain or lose from an interaction, we need to enumerate the possible ways in which the trustee may act, be it to the advantage or disadvantage of the truster. In addition, the granularity with which we represent a trustee's actions depends on what matters to the truster.

For instance, if Bob has 50 pence to spend on a phone call to a long lost friend, then the more time he can spend on the phone for that money, the better it may be for him. On the other hand, suppose that Bob is just calling to order a pizza, and that, regardless of how long he takes to place his order, he will not receive any change from his 50 pence. In this case, perhaps it doesn't matter if his credit runs out in 5 minutes or 5 hours, provided he has enough time to place his order.

This shows that sometimes an agent's preferences may depend on one or more finely graded attributes of a trustee's behaviour, or it may depend only on whether an action is carried out or not. In the latter case, it is sufficient to represent a trustee's actions as a binary event: either the trustee cooperates by fulfilling its obligations to the truster, or it defects by breaking its obligations.

With this in mind, we give separate treatments for each of the two types of cases described, by concentrating on the binary case in the current chapter, and leaving the continuous case to Chapter 5. Specifically, the current chapter introduces a trust model that we call TRAVOS (Trust and Reputation system for Agent Based Virtual OrganisationS), which instantiates the framework described in the previous chapter, for boolean action spaces.

To begin our discussion, Sections 4.1 and 4.2 outline how, by applying the framework, Bayesian analysis is used in TRAVOS to assess a trustee based on direct experience and reputation. Using this as a foundation, we then describe the main contributions of TRAVOS to the state of the art, which comes in three parts. First, Section 4.3 describes how inaccurate reputation is handled by TRAVOS. To achieve this, a truster learns over time which of its reputation sources are reliable, by comparing the opinions provided about a trustee to its subsequent behaviour. Based on this, if a reputation source generally provides opinions that do not correlate well with actual trustee behaviour, it will have its impact on the truster’s future assessments reduced. Second, Sections 4.4 and 4.5 describe how TRAVOS can be used as part of the CONOISE-G¹ system, for forming and managing virtual organisations in a service-oriented environment. Third, Section 4.6 demonstrates empirically how TRAVOS outperforms the most similar existing model in the literature.

4.1 Instantiating the Framework for Boolean Action Spaces

The process of applying the principles of the previous chapter to a particular problem consists of three main steps: (1) we need to choose an appropriate parameter model for representing interaction outcomes, along with their distributions; (2) we need to decide how to communicate reputation between agents, based on the guidelines laid down previously; and (3) we need to specify a mechanism for dealing with inaccurate reputation sources.

Thus, we start by describing the parameter model used in TRAVOS to reason about interaction outcomes that, as we have already stated, can take on one of two values: a trustee can either cooperate by fulfilling its obligations, or it can defect by neglecting its obligations. As such, this description of behaviour can naturally be represented by a binary number, which we realise in our notation by setting $\mathcal{O}^c = \{0, 1\}$, and attaching the following semantics to any $O_{a_{tr}, a_{te}} \in \mathcal{O}^c$.

$$O_{a_{tr}, a_{te}} = \begin{cases} 1 & \text{if contract is fulfilled by } a_{te} \\ 0 & \text{otherwise} \end{cases} \quad (4.1)$$

This binary definition means that a series of observations of trustee behaviour (such as $O_{a_{tr}, a_{te}}^{0:t'}$) can be treated in the same way as a series of tosses from a biased coin — that is, as a set of Bernoulli trials. We have already described an appropriate parameter model for this in Chapter 3, in which the distribution of outcomes is described by a single parameter representing the probability of obtaining a 1 or 0. Thus, we define the behaviour distribution parameter, $\theta_{a_{tr}, a_{te}}$, as the probability that a_{te} will fulfill its

¹The name CONOISE-G stands for Constraint Oriented Negotiation in Open Information Seeking Environments for the Grid.

obligations during an interaction with a_{tr} (Equation 4.2), and specify its domain as $\Theta^C = [0, 1]$:

$$\theta_{a_{tr}, a_{te}} = p(O_{a_{tr}, a_{te}} = 1), \quad \text{where } \theta_{a_{tr}, a_{te}} \in \Theta^C = [0, 1] \quad (4.2)$$

In the interest of simplicity, we adopt the standard practice of choosing a conjugate prior for the parameter distribution (DeGroot and Schervish, 2002) that, in the case of Bernoulli distributions, is given by the family of beta distributions. In this respect, TRAVOS is therefore similar to the Beta Reputation System and its related models (Section 2.2.2), which also represent interaction outcomes as binary events, and model their parameter distributions as Beta distributions.

Thus, the hyperparameter space, Φ^C , now takes on the form of the standard parameters of the beta distribution (Equation 4.3). Specifically, the beta distribution has two parameters, typically denoted α and β , both of which are positive real numbers. These parameters determine the shape of the distribution through the probability density function (Equation 4.4), the expected value of the distribution (Equation 4.5) and the variance (Equation 4.6).

$$\Phi^C = \{(\alpha, \beta) | \alpha > 0 \wedge \beta > 0\} \quad (4.3)$$

$$d(\theta | \alpha, \beta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int U^{\alpha-1} (1-U)^{\beta-1} dU} \quad (4.4)$$

$$E[\theta | \alpha, \beta] = \frac{\alpha}{\alpha + \beta} \quad (4.5)$$

$$\sigma^2 = \frac{\alpha \cdot \beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)} \quad (4.6)$$

With this in mind, we can now show how the various aspects of the beta distribution can be applied to the framework. In particular, for a given prior, $\phi_{a_{tr}, a_{te}}^{prior} = (\alpha^{prior}, \beta^{prior})$, the posterior hyperparameter, $\phi_{a_{tr}, a_{te}}^{post} = (\alpha^{post}, \beta^{post})$, is calculated by counting the number of successful interactions (Equation 4.7) and the number of unsuccessful interactions (Equation 4.8) in the interaction history, $O_{a_{tr}, a_{te}}^{0:t'}$, and then adding these values to the α and β parameters as shown in Equations 4.9 and 4.10. This is a well known result, a derivation of which is given by DeGroot and Schervish (2002).

$$m_{a_{tr}, a_{te}} = |\{o \in O_{a_{tr}, a_{te}}^{0:t'} | o = 1\}| \quad (4.7)$$

$$n_{a_{tr}, a_{te}} = |\{o \in O_{a_{tr}, a_{te}}^{0:t'} | o = 0\}| \quad (4.8)$$

$$\alpha^{post} = \alpha^{prior} + m_{a_{tr}, a_{te}} \quad (4.9)$$

$$\beta^{post} = \beta^{prior} + n_{a_{tr}, a_{te}} \quad (4.10)$$

The effect of updating the parameter distribution in light of observations is illustrated in Figure 4.1. Here, adding observations, and thus increasing α and β , decreases the distribution variance, making the distribution more peaked. The proportion of successful and unsuccessful interactions, along with the prior, determine where in the interval $[0, 1]$

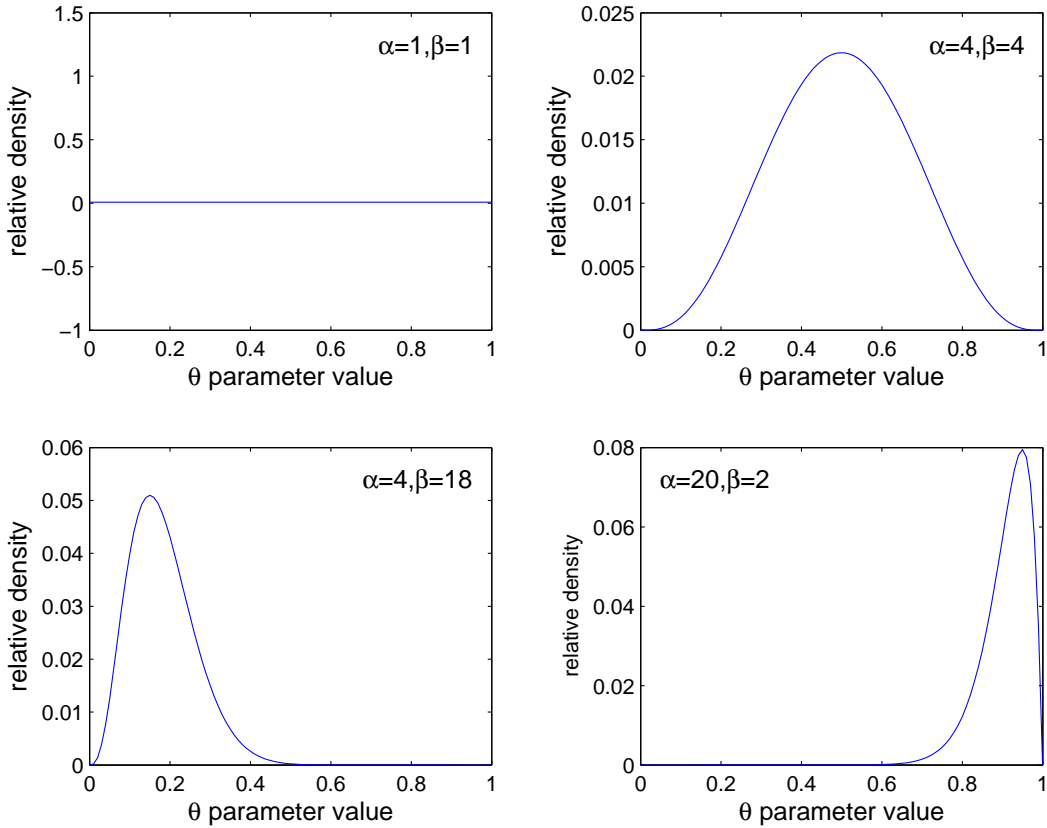


FIGURE 4.1: Example beta pdf plots; note that when $\alpha = 1, \beta = 1$ (top-left) the distribution is uniform in the interval $[0, 1]$.

the distribution peaks. A high α value compared to β (usually resulting from a high proportion of successful outcomes) causes the distribution mode to occur close to 1. Intuitively, this is correct because it supports the conclusion that the intrinsic probability of $O_{a_{tr}, a_{te}} = 1$ is also close to 1.

At this stage, it is useful to describe in detail how the parameter model in TRAVOS can be applied within the framework of Decision Theory. Section 3.3 described this in more general terms by stipulating that a truster's preferences should be captured in a utility function, $U : \mathcal{O}^C \rightarrow \mathbb{R}$, and that this should be used, along with the distribution of $O_{a_{tr}, a_{te}}$, to calculate the expected utility of a_{tr} interacting with a_{te} . Assuming that $p(\theta_{a_{tr}, a_{te}})$ is the p.d.f. of $\theta_{a_{tr}, a_{te}}$ given all available direct experience and reputation, then it follows from the parameter model that this expected utility can be calculated in the

following way:

$$EU = \sum_{o \in \mathcal{O}^c} U(O_{a_{tr}, a_{te}} = o) p(O_{a_{tr}, a_{te}} = o) \quad (4.11)$$

$$= \sum_{o \in \mathcal{O}^c} U(O_{a_{tr}, a_{te}} = o) \int_{\Theta^c} p(O_{a_{tr}, a_{te}} = o | \theta_{a_{tr}, a_{te}}) p(\theta_{a_{tr}, a_{te}}) d\theta_{a_{tr}, a_{te}} \quad (4.12)$$

$$= U(O_{a_{tr}, a_{te}} = 1) \int_{\Theta^c} \theta_{a_{tr}, a_{te}} p(\theta_{a_{tr}, a_{te}}) d\theta_{a_{tr}, a_{te}} + \quad (4.13)$$

$$U(O_{a_{tr}, a_{te}} = 0) \int_{\Theta^c} (1 - \theta_{a_{tr}, a_{te}}) p(\theta_{a_{tr}, a_{te}}) d\theta_{a_{tr}, a_{te}}$$

$$= E[\theta_{a_{tr}, a_{te}}] U(O_{a_{tr}, a_{te}} = 1) + (1 - E[\theta_{a_{tr}, a_{te}}]) U(O_{a_{tr}, a_{te}} = 0) \quad (4.14)$$

What is interesting about this calculation is its simplicity, and that it depends only on the expected value of $\theta_{a_{tr}, a_{te}}$ (Equation 4.5), rather than the parameter distribution as a whole. Choosing between competing trustees thus becomes the simple process of opting for the agent that maximises Equation 4.14.

4.2 Instantiating the Reputation Mechanism

Accounting for reputation in TRAVOS requires agents to be able to share their experiences of one another's behaviour. According to our guidelines, this should be achieved by specifying an opinion function $R_{a_{rep}, a_{te}} = r(O_{a_{rep}, a_{te}}^{0:t'})$ that has a shared definition across all reputation sources $a_{rep} \in \mathcal{A}$, and is (ideally) a minimal sufficient statistic of $O_{a_{rep}, a_{te}}^{0:t'}$.

To achieve this in TRAVOS, we define $R_{a_{rep}, a_{te}}$ to be a two dimensional vector consisting of the number of successful and unsuccessful interactions that a_{rep} has had with a_{te} , or put another way, we let $R_{a_{rep}, a_{te}}$ be the vector $\langle m_{a_{rep}, a_{te}}, n_{a_{rep}, a_{te}} \rangle$, in which each component is defined as follows:

$$m_{a_{rep}, a_{te}} = |\{o \in O_{a_{rep}, a_{te}}^{0:t'} | o = 1\}| \quad (4.15)$$

$$n_{a_{rep}, a_{te}} = |\{o \in O_{a_{rep}, a_{te}}^{0:t'} | o = 0\}| \quad (4.16)$$

Accounting for reputation in our assessments can now be done in the same way as a truster does for its own direct experiences. That is, we simply sum together the truster's own interaction counts with those of its reputation sources, and so calculate the posterior

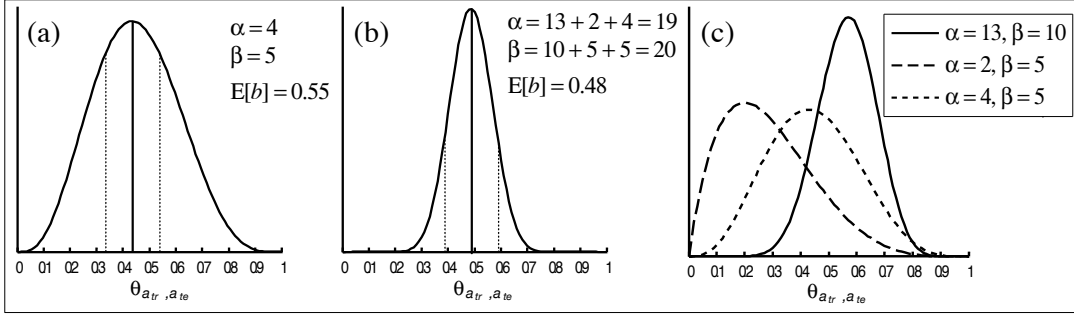


FIGURE 4.2: Example beta distributions for aggregating opinions of 3 agents.

hyperparameters in the following way, given reputation sources $\{a_i, \dots, a_p\}$.

$$\alpha^{post} = M_{a_{tr}, a_{te}} + \alpha^{prior} \quad (4.17)$$

$$\beta^{post} = N_{a_{tr}, a_{te}} + \beta^{prior} \quad \text{where,} \quad (4.18)$$

$$M_{a_{tr}, a_{te}} = m_{a_{tr}, a_{te}} + \sum_{i=1}^p m_{a_i, a_{te}} \quad (4.19)$$

$$N_{a_{tr}, a_{te}} = n_{a_{tr}, a_{te}} + \sum_{i=1}^p n_{a_i, a_{te}} \quad (4.20)$$

From this, it is clear that the opinion function is sufficient because the posterior distribution depends only on the interaction counts, and not on any other properties of the observations. In addition, it is minimal because, with all other factors constant, any change in a reputation source's opinion will always result in a change to the posterior hyperparameters. As such, there is a one-to-one mapping between an opinion and its effect on the hyperparameters, so any non-invertible function of $R_{a_{rep}, a_{te}}$ would not be sufficient, thus proving that $R_{a_{rep}, a_{te}}$ is minimal.

The effect of combining opinions in this way is illustrated in Figure 4.2. In this figure, part (a) shows a beta distribution representing one agent's opinion, along with the attributes of the distribution that have been discussed so far. In contrast to this, part (c) illustrates the differences between the distribution in part (a) and distributions representing the opinions of two other agents with different experiences. The result of combining all three opinions is illustrated in part (b), of which there are two important characteristics. First, the distribution with parameters $\alpha = 13$ and $\beta = 10$ is based on more observations than the remaining two distributions put together, and so has the greatest impact on the shape and expected value of the combined distribution. This demonstrates how conflicts between different opinions are resolved: the combined trust value is essentially a weighted average of the individual opinions, where opinions with higher confidence values are given greater weight. Second, the variance of the combined distribution is strictly less than any one of the component distributions. This reflects the fact that it is based on more observations overall, and so has a greater confidence value.

4.3 Filtering Inaccurate Reputation

As described in the previous chapter, opinions provided about a trustee from a third party may not always be reliable. This may occur for a variety of reasons, for example, because a trustee behaves differently toward a third party than it does toward the truster, or because a reputation source intentionally manipulates its opinions for its own purposes. In addition, the latter of these possibilities is particularly problematic because not only do we not know for certain that reported observations are drawn from the desired distribution, but also because the number of observations may be exaggerated to increase their impact on a truster's assessment.

In Chapter 2, we reviewed two basic approaches for addressing this problem, which Jøsang et al. (2005) refer to as *endogenous* and *exogenous* methods. The former attempt to identify unreliable reputation information by considering the statistical properties of the reported opinions alone (e.g. Whitby et al. (2004); Dellarocas (2000)), while the latter rely on other information to make such judgements, such as the reputation of the source or its relationship with the trustee (e.g. Buchegger and Boudec (2003); Yu and Singh (2003); Klos and Poutré (2004)).

Many proposals for endogenous techniques assume that inaccurate or unfair raters are generally in a minority among reputation sources, and thus consider reputation providers whose opinions deviate in some way from mainstream opinion to be those most likely to be inaccurate. Our solution is exogenous, in that we judge a reputation provider on the perceived accuracy of its past opinions, rather than its deviation from mainstream opinion. Moreover, we define a two-step method as follows. First, we calculate the probability that an agent will provide an accurate opinion given its past opinions and later observed² interactions with the trustees for which opinions were given. Second, based on this value, we reduce the distance between a rater's opinion and the prior belief that all possible values for an agent's behaviour are equally probable. Once all the opinions collected about a trustee have been adjusted in this way, the opinions are aggregated using the technique described above. In so doing, we reduce the influence that an opinion provider has on a truster's assessment of a trustee, if the provider's opinion is consistently biased in one way or another. This can be true either if the provider is malevolent, or if a significant number of trustees behave differently towards the truster than towards the opinion provider in question.

We describe this technique in more detail in the remainder of this section: first we detail how the probability of accuracy is calculated, and then we show how opinions are adjusted and the combined reputation obtained. An example of how these techniques can be used is also given with the aid of a walk-through scenario in Section 4.5.2.

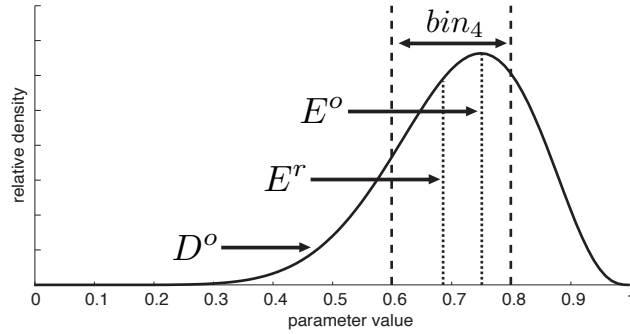
²These are observations made by the truster after it has obtained an opinion.

4.3.1 Estimating the Probability of Accuracy

The first stage in our solution is to estimate the probability that a rater's stated opinion of a trustee is accurate, which depends on the value of the current opinion under consideration, denoted $\hat{\mathcal{R}}_{a_{rep}, a_{te}} = (\hat{m}_{a_{rep}, a_{te}}, \hat{n}_{a_{rep}, a_{te}})$. Specifically, if E^r is the expected value of a beta distribution, D^r , such that $\alpha^r = \hat{m}_{a_{rep}, a_{te}} + 1$ and $\beta^r = \hat{n}_{a_{rep}, a_{te}} + 1$, we can estimate the probability that E^r lies within some margin of error around $\theta_{a_{tr}, a_{te}}$, which we call the accuracy of a_{rep} according to a_{tr} , denoted as $\rho_{a_{tr}, a_{rep}}$. To perform this estimation, we consider the outcomes of all previous interactions for which a_{rep} provided an opinion similar to $\hat{\mathcal{R}}_{a_{rep}, a_{te}}$ about a_{te} , to a_{tr} , for each a_{te} . Using these outcomes, we construct a beta distribution, D^o , for which, if its expected value E^o is close to E^r , then a_{rep} 's opinions are generally correlated to what is actually observed, and we can judge a_{rep} 's accuracy to be high. Conversely, if E^r deviates significantly from E^o , then a_{rep} has low accuracy.

For example, the process of achieving this estimation is illustrated in Figure 4.3, in which the range of possible values of E^r and E^o is divided into five intervals (or bins), $bin_1 = [0, 0.2], \dots, bin_5 = [0.8, 1]$. These bins define which opinions we consider to be similar to each other, such that all opinions that lie in the same bin are considered alike. This is necessary because we may never see enough opinions from the same provider to assess an opinion based on identical opinions in the past. Instead, the best we can do is consider the perceived accuracy of past opinions that do not deviate significantly from the opinion under consideration. In the case illustrated in the figure, the opinion provider, a_{rep} , has provided a_{tr} with an opinion with an expected value in bin_4 . Now, if we therefore consider all previous interaction outcomes for which a_{rep} provided an opinion to a_{tr} in bin_4 , the portion of successful outcomes, and thus E^o , is also in bin_4 , so $\rho_{a_{tr}, a_{rep}}$ is high. If subsequent outcome-opinion pairs were also to follow this trend, then D^o would be highly peaked inside this interval, and $\rho_{a_{tr}, a_{rep}}$ would converge to 1. Conversely, if subsequent outcomes disagreed with their corresponding opinions, then $\rho_{a_{tr}, a_{rep}}$ would approach 0.

More specifically, we divide the range of possible values of E^r into N disjoint intervals bin_1, \dots, bin_n , then calculate E^r , and find the interval, bin^o , that contains the value of E^r . Then, if $\mathcal{H}_{a_{tr}, a_{rep}}$ is the set of all pairs of the form $(O_{a_{tr}, a_x}, \hat{\mathcal{R}}_{a_{rep}, a_x})$, where $a_x \in \mathcal{A}$, and O_{a_{tr}, a_x} is the outcome of an interaction for which, prior to being observed by a_{tr} , a_{rep} gave the opinion $\hat{\mathcal{R}}_{a_{rep}, a_x}$, we can find the subset $\mathcal{H}_{a_{tr}, a_{rep}}^r \subseteq \mathcal{H}_{a_{tr}, a_{rep}}$, which comprises all pairs for which the opinion's expected value falls in bin^o . We then count the total number of pairs in $\mathcal{H}_{a_{tr}, a_{rep}}^r$ for which the interaction outcome was successful (denoted $C_{success}$) and those for which it was not (denoted C_{fail}). Based on these frequencies, the parameters for D^o can be defined as $\alpha^o = C_{success} + 1$ and $\beta^o = C_{fail} + 1$. Using D^o , we now calculate $\rho_{a_{tr}, a_{rep}}$ as the portion of the total mass of D^o that lies in the interval

FIGURE 4.3: Illustration of $\rho_{a_{tr}, a_{rep}}$ estimation process.

bin^o (see Equation 4.21).

$$\rho_{a_{tr}, a_{rep}} = \frac{\int_{\min(bin^o)}^{\max(bin^o)} X^{\alpha^o-1} (1-X)^{\beta^o-1} dX}{\int_0^1 U^{\alpha^o-1} (1-U)^{\beta^o-1} dU} \quad (4.21)$$

Each trustor performs these operations to determine the probability of accuracy of reported opinions. However, one implication of this technique is that the number (and size) of bins effectively determines an acceptable margin of error in opinion provider accuracy: the estimated accuracy of a larger set of opinion providers converges to 1 with large bin sizes, as opposed to small sizes.

4.3.2 Adjusting Reputation Source Opinions

To describe how we adjust reputation opinions, we must introduce some new notation. First, let D^c be the beta distribution that results from combining all of a trustee's reputation information (using Equations 4.17 to 4.20). Second, let D^{c-r} be a distribution constructed using the same equations, except that the opinion under consideration, $\hat{\mathcal{R}}_{a_{rep}, a_{te}}$, is omitted. Third, let \bar{D} be the result of adjusting the opinion distribution D^r , according to the process described below. Finally, we refer to the standard deviation (denoted σ), expected value and parameters of each distribution by using the respective superscript; for instance, D^c has parameters α^c and β^c , with standard deviation σ^c and expected value E^c .

Now, our goal is to reduce the *effect* of unreliable opinions on D^c . In essence, by adding $\hat{\mathcal{R}}_{a_{rep}, a_{te}}$ to a trustee's reputation, we move E^c in the direction of E^r . The standard deviation of D^r contributes to the confidence value for the combined reputation value but, more importantly, its value relative to σ^{c-r} determines how far E^c will move towards E^r . This effect has important implications: consider as an example three distributions d_1 , d_2 and d_3 , with shape parameters, expected value and standard deviation as shown in Table 4.1; the results of combining d_1 with each of the other two distributions are

Distribution	α	β	E	σ
d_1	540	280	0.6585	0.0165
d_2	200	200	0.5000	0.0250
d_3	5000	5000	0.5000	0.0050
$d_1 + d_2$	740	480	0.6066	0.0140
$d_1 + d_3$	5540	5280	0.5120	0.0048

TABLE 4.1: Combination of beta distributions.

shown in the last two rows. As can be seen, distributions d_2 and d_3 have identical expected values with standard deviations of 0.025 and 0.005 respectively. Although the difference between these values is small (0.02), the result of combining d_1 with d_2 is quite different from combining d_1 and d_3 . Whereas the expected value in the first case falls approximately between the expected values for d_1 and d_2 , the relatively small parameter values of d_1 compared to d_3 in the latter case means that d_1 has virtually no impact on the combined result. Obviously, this is due to our method of reputation combination in which the parameter values are summed. This is important because it shows how, if left unchecked, an unfair rater could deliberately increase the weight an agent places on its opinion by providing very large values for m and n which, in turn, determine α and β .

In light of this, we adopt an approach that significantly reduces very high parameter values unless the probability of the rater’s opinion being accurate is very close to 1. Specifically, we reduce the distance between, respectively, the expected value and standard deviation of D^r , and the expected value and standard deviation of the uniform distribution, $\alpha = \beta = 1$, which represents a state of no information (see Equations 4.22 and 4.23). Here, we denote the standard deviation of the uniform distribution as $\sigma_{uniform}$ and its expected value as $E_{uniform}$. By adjusting the standard deviation in this way, rather than changing the α and β parameters directly, we ensure that large parameter values are decreased more than smaller values. We adjust the expected value to guard against cases where we do not have enough reliable opinions to mediate the effect of unreliable opinions; if we did not adjust the expected value then, in the absence of any other information, we would take an opinion source’s word as true, even if we did not consider its opinion reliable.

$$\bar{E} = E_{uniform} + \rho_{a_{tr}, a_{rep}} \cdot (E^r - E_{uniform}) \quad (4.22)$$

$$\bar{\sigma} = \sigma_{uniform} + \rho_{a_{tr}, a_{rep}} \cdot (\sigma^r - \sigma_{uniform}) \quad (4.23)$$

Once we have determined the values of \bar{E} and $\bar{\sigma}$, we use Equations 4.24 and 4.25 to find the parameters $\bar{\alpha}$ and $\bar{\beta}$ of the adjusted distribution,³ and from these we calculate adjusted values for $\hat{m}_{a_{rep}, a_{te}}$ and $\hat{n}_{a_{rep}, a_{te}}$, denoted as $\bar{m}_{a_{rep}, a_{te}}$ and $\bar{n}_{a_{rep}, a_{te}}$ respectively (see Equation 4.26). These scaled versions of $\hat{m}_{a_{rep}, a_{te}}$ and $\hat{n}_{a_{rep}, a_{te}}$ are then used in their place to calculate the combined trust value, as in Equations 4.17 to 4.20. Strictly

³A derivation of these equations is provided in Appendix B.

speaking, $\bar{m}_{a_{rep},a_{te}}$ and $\bar{n}_{a_{rep},a_{te}}$ are not frequencies nor are their unadjusted counterparts, but, for these purposes, these have the same effect on the combined trust value as an equivalent set of observations made by the truster itself. In general, as $\rho_{a_{tr},a_{rep}}$ approaches 0, both $\bar{m}_{a_{rep},a_{te}}$ and $\bar{n}_{a_{rep},a_{te}}$ will also approach 0. Thus, if $\rho_{a_{tr},a_{rep}}$ is 0 then no observation reported by a_{rep} will affect a_{tr} 's decision making in any way.

$$\bar{\alpha} = \frac{\bar{E}^2 - \bar{E}^3}{\bar{\sigma}^2} - \bar{E} \quad (4.24)$$

$$\bar{\beta} = \frac{(1 - \bar{E})^2 - (1 - \bar{E})^3}{\bar{\sigma}^2} - (1 - \bar{E}) \quad (4.25)$$

$$\bar{m}_{a_{rep},a_{te}} = \bar{\alpha} - 1 \quad , \quad \bar{n}_{a_{rep},a_{te}} = \bar{\beta} - 1 \quad (4.26)$$

4.4 Reputation Gathering for TRAVOS

In the preceding sections, we have shown how, by using the framework, reputation information can be used along with a truster's direct experience to assess the trustworthiness of an agent. However, apart from assessment, there are two other issues that a practical trust and reputation system should include: (1) agents require some mechanism to obtain opinions from reputation sources, and (2) agents must decide when it is necessary to obtain reputation information. The latter is important, because if a truster has sufficient direct evidence with which to judge a trustee, the cost of obtaining reputation information may outweigh its benefits. We now consider each of these issues in turn.

4.4.1 Reputation Brokering

The problem with obtaining opinions in large systems is that directly querying many agents may entail a significant communication overhead. Therefore, agents must do one or more of the following:

1. choose a subset of agents to query; and
2. employ some method of streamlining reputation.

Our solution to this problem is illustrated in Figure 4.4. We assume that each agent in a system belongs to exactly one primary *domain*. Here, a domain may correspond to an organisation or department in the real world, to which the agent is responsible. This view is in line with the vision of systems, such as the Grid, in which computing resources belonging to different organisations may be used together (recall our discussion in Section 1.3). Within each domain, there is a *reputation broker* agent, which is responsible for aggregating the opinions of all other agents within its domain; that is, the opinion of a reputation broker about a trustee is an aggregation of the opinions of

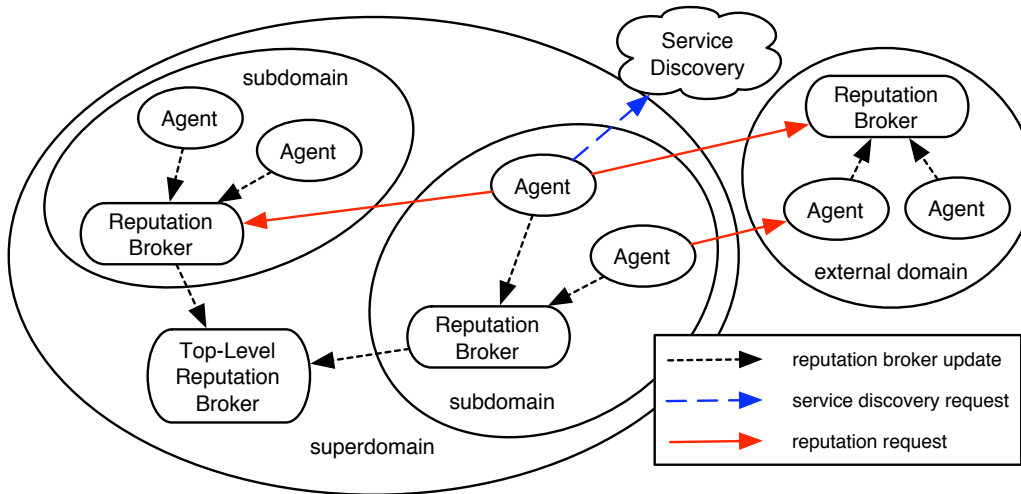


FIGURE 4.4: Reputation brokering system.

all other agents within its domain. In addition, domains can be arranged in a hierarchy such that brokers in sub-domains report to a broker in an overall domain. In this way, a top level broker aggregates all the opinions of agents in each of its sub-domains.

Reputation brokers provide a point of contact for external agents looking to receive reputation information. When a trustor requires reputation, it first uses a service discovery system (such as described in Section 4.5) to identify domains that advertise information about trustees in some general context⁴. For example, companies which make use of grid-based storage space may advertise knowledge about vendors of such storage space.

Once a trustor has received a list of appropriate domains, it can choose to request an opinion from either the main reputation broker for that domain, or other brokers or individual agents within that domain. Although we do not specify how a trustor should make this choice, there is an obvious trade-off in granularity. By requesting information from a top-level broker, the trustor can receive all the information known by the domain in a single message. However, in this case, a trustor can only judge the accuracy of the broker's domain as a whole (using the techniques described in Section 4.2). On the other hand, if a trustor contacted several agents within a domain, it could judge their accuracy individually, thus identifying the most reliable contacts within an organisation. Here, it is important that a trustor should avoid using a reputation source at the same time as any reputation broker that the source reports to. The reason for this is the problem of correlated evidence: since the broker's opinion is based on those agents which report to it, using a reputation source along with its broker would amount to counting the reputation source's opinion twice!

⁴Here, we do not address the issue of the level at which domains should advertise information. For example, if a department within a company is mainly responsible for certain information, it is not clear whether the department should be the advertised point of contact, or the organisation it belongs to.

Algorithm 1 Reputation broker update algorithm, performed by reputation sources.

```

{Each time an interaction outcome is observed do the following}
if interaction successful then
   $m[\text{trustee\_id}] \leftarrow m + 1$ 
else
   $n[\text{trustee\_id}] \leftarrow n + 1$ 
end if
{Periodically do the following}
for all  $i = \text{trustee\_id}$  do
  if  $m[i] \neq 0$  OR  $N[I] \neq 0$  then
    add  $m[i]$  and  $n[i]$  to update message
  end if
   $m[i] \leftarrow 0$ 
   $n[i] \leftarrow 0$ 
end for
SEND update message to reputation broker

```

We now describe how a reputation broker's opinion is formed. Each broker periodically receives updates regarding any newly observed interaction outcomes from the agents within its own domain of responsibility. These updates take the same form as normal reputation opinions in TRAVOS (Equations 4.17 to 4.20) except that they are only based on observations that have occurred since the last update the observer sent to its broker. This process is summarised in Algorithm 1. When a reputation broker receives an opinion from within its domain about a trustee a_{te} , it updates its own opinion about a_{te} using Equations 4.27 and 4.28. In this way, the broker's opinion can be compared to that of a single agent that has observed all the interaction outcomes recorded by the agents within the broker's domain.

$$m_{a_{broker}, a_{te}} = m_{a_{broker}, a_{te}} + \sum_{a_i \in \mathcal{D}} m_{a_i, a_{te}}^* \quad (4.27)$$

$$n_{a_{broker}, a_{te}} = n_{a_{broker}, a_{te}} + \sum_{a_i \in \mathcal{D}} n_{a_i, a_{te}}^* \quad (4.28)$$

where $(m_{a_i, a_{te}}^*, n_{a_i, a_{te}}^*)$ is the update message from a_i about a_{te} ,

and \mathcal{D} is the set of agents in a_{broker} 's domain.

4.4.2 When to Seek Reputation

In some cases, an agent may decide that it is sufficiently confident in its own knowledge about a trustee to avoid acquiring reputation information to improve its estimate. This is valuable because of the communication cost of reputation acquisition, and the inherent unreliability of reputation compared to direct observations. One simple method of doing this is to calculate the posterior probability that the true value for $\theta_{a_{tr}, a_{te}}$ lies within an acceptable margin of error around its estimate, $\vartheta_{a_{tr}, a_{te}}$. We can calculate this using the parameter distribution as follows. First, we decide on an acceptable error margin,

$\vartheta_{a_{tr}a_{te}} \pm \epsilon$, where ϵ is the acceptable distance from $\vartheta_{a_{tr}a_{te}}$. Second, we integrate the parameter distribution over the area defined by the error margin. To do this, we use the beta probability density function as shown in Equation 4.29. We refer to the resulting value as the *confidence* value of $\vartheta_{a_{tr}a_{te}}$, which we denote as $\gamma_{a_{tr},a_{te}}$. Finally, we choose a threshold for this probability, above which we consider the accuracy of the estimate as acceptable; we denote this threshold as τ .

$$\gamma_{a_{tr},a_{te}} = \frac{\int_{\vartheta_{a_{tr},a_{te}} - \epsilon}^{\vartheta_{a_{tr},a_{te}} + \epsilon} B^{\alpha-1}(1-B)^{\beta-1} dB}{\int_0^1 U^{\alpha-1}(1-U)^{\beta-1} dU}, \quad \text{where } (\alpha, \beta) = \phi_{a_{tr},a_{te}} \quad (4.29)$$

4.5 An Application to Agent-Based Virtual Organisations

In this section, we describe the role of TRAVOS in the CONOISE-G system (Patel et al., 2005b; Shao et al., 2004), which seeks to “*support robust and resilient virtual organisation formation and operation. It aims to provide mechanisms to assure effective operation of VOs in the face of disruptive and potentially malicious entities in dynamic, open and competitive environments.*”⁵ More specifically, CONOISE-G provides methods by which agents operating in a grid environment can form *dynamic* resource coalitions (VOs) in order to fulfill their goals. Here, by dynamic we mean that the membership of a VO may change over its lifetime. This can happen for various reasons; for instance, a particular member’s resources may fail, requiring a new member to make up the shortfall. In the following subsections we give an overview of the CONOISE-G system (Section 4.5.1), followed by a trust-oriented scenario of how TRAVOS is used in CONOISE-G (Section 4.5.2). All of the work described in these sections has been implemented, and applied to a number of real world scenarios.

4.5.1 System Overview

In essence, the CONOISE-G architecture comprises several different agents, including *system agents* and *service providers* (SPs), as shown in Figure 4.5. SPs are agents which may belong to a VO, and are responsible for overseeing its formation, operation and dissolution. On the other hand, system agents are those needed to achieve core system functionality for VO formation and operation, and include five different types of agent:

1. the Yellow Pages Agent (YP);
2. the Quality of Service Consultant (QoS);
3. the Quality of Service Accessor (QA);

⁵This quote is taken from <http://www.conoise.org/>

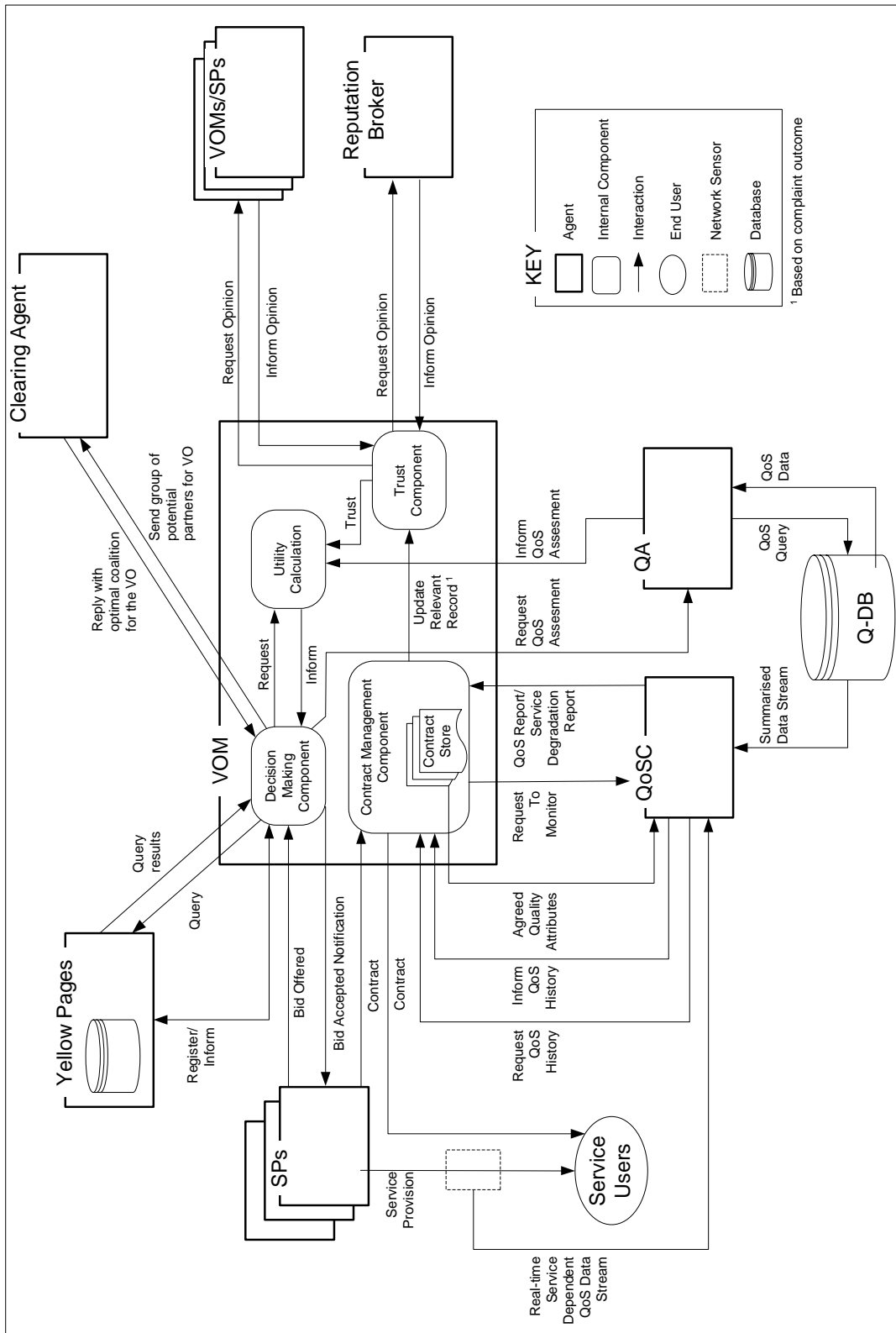


FIGURE 4.5: The CONOISE-G architecture.

4. the Clearing Agent (CA); and
5. Reputation Brokers.

Using these agents, the formation of a VO consists of three steps:

1. **Resource Discovery** — A particular SP, acting either on its own behalf, or on behalf of a user, identifies a need for a number of resources, which it cannot supply (efficiently) by itself. To fulfill this need, the SP instigates VO formation, by requesting a list of other SPs that can supply the required resources; it obtains this list from the Yellow Pages Agent (YP), which performs a service discovery role (Deora et al., 2004). At this point, the SP that places the request for the service takes on the VO Manager (VOM) role for the potential VO, as illustrated in Figure 4.5.
2. **Resource Assessment** – After receiving a response from the YP, the VOM invites the identified providers to bid for the requested services. Once all such bids are received, the VOM generates an expected utility function for each bid, based on the price offered per resource unit, trust and the advice given by the Quality of Service Assessor (QoSA). The QoSA, based on Deora et al. (2003), is an external service which rates how well a given SP is likely to perform. Its role can be viewed as similar to that of a reputation provider in TRAVOS, in that it provides extra information about a trustee’s likely behaviour. However, the nature of its assessment and its underlying assumptions are different from that of reputation sharing in TRAVOS, and therefore it must be treated differently.

In our approach, we first estimate the SP’s behaviour distribution (as described in the previous sections) thereby estimating the probability that the SP will fulfill its obligations to the VOM. Then, we use the QoSA’s assessment of an SP to provide an alternative estimate of this probability, and combine these two estimates using a suprabayesian approach (Keeney and Raiffa, 1976). In general, the combined probability should be more accurate than either of the individual estimates, since it incorporates the knowledge of the QoSA, the VOM (in its role as a trustor) and the VOM’s reputation sources. The combined probability is then used to calculate the expected utility for the VOM, for each possible number of resource units it can purchase from the bidding SP.

3. **Resource Allocation** – Once we have an expected utility function for each bidder, we employ the Clearing Agent (CA), which finds the optimal resource allocation⁶ for the set of bidding SPs (Dang and Jennings, 2002). The resulting allocations are reported back to the VOM, which then sends ‘hired’ messages to

⁶Alternatively, if there are significant time constraints, the CA can find an allocation which is within some bound of optimal.

each of the successful bidders, informing them of the quantity of each resource they are asked to provide.

Once the VO is formed, the operational phase begins. During this stage, the VOM may request the QoS Consultant (QoSC) to monitor any services provided by any members of the VO. The QoSC informs the VOM if and when an SP diverges from its agreed service level. When the QoS provision of a service in the VO falls below an acceptable level of service, or some breach of contract is observed, the QoSC alerts the VOM, which initiates a VO re-formation process. During this stage, the Contract Management component of the VOM decides whether a breach of contract has actually occurred, and if so, which SP is to blame. Based on this result, the VOM updates its trust component, recording either a successful or unsuccessful outcome for any terminated contracts.

Meanwhile, the VOM issues another message to the YP, requesting a list of SPs that can replace the resources of the failed SP. As before, the YP identifies possible SPs, bids are received and evaluated, resulting in the CA determining the best SP to replace the failed provider. At this point, the VOM re-forms the VO with the new SP replacing the old one, and instructs the QoSC to stop monitoring the old SP and to monitor the new one instead. A similar process may also take place if another SP, not currently in the VO, sends the VOM a competitive offer on resources it receives from current VO members. This process is facilitated by a publish and subscribe service offered by the YP: the VOM may register interest in SPs that provide particular resources, in response to which the YP will inform the VOM any time a new SP offering such services appears in the system.

4.5.2 Walk-through Scenario

This section provides an agent-based VO scenario in which we demonstrate the use of TRAVOS. We begin by stating that there is a need to create a VO to meet a specific requirement to provide a composite multimedia communication service to an end user. This consists of the following basic services: text messaging, multimedia streaming, HTML content provision, and phone calls (this example is adapted from one given by [Norman et al. \(2003\)](#)). Now, assume agent a_1 has identified this need and wishes to capitalise on the market niche. However, a_1 only has the capability to provide a text messaging service, and can only achieve its goal by forming a VO with an agent that can supply a service for phone calls and one for HTML content. For simplicity, we assume that each agent in the system has the ability to provide only one service. Agent a_1 is aware of three agents that can provide a phone call service, and its interaction history with these is shown in [Table 4.2](#). Similarly, it is aware of three agents that are capable of providing HTML content, and its past interactions with these entities are given in [Table 4.3](#). We also assume that a trustor's prior parameter distribution for all agents is

Agent	Past interactions	
	Successful	Unsuccessful
a_2	17	5
a_3	2	15
a_4	18	5

TABLE 4.2: Agent a_1 's interaction history with phone call service provider agents.

Agent	Past interactions	
	Successful	Unsuccessful
a_5	9	14
a_6	3	0
a_7	18	11

TABLE 4.3: Agent a_1 's interaction history with HTML content service provider agents.

uniform:

$$\alpha^{prior} = 1, \quad \beta^{prior} = 1$$

Agent a_1 would like to choose the most trustworthy phone call and HTML content service provider from the selection. The following describes how this is achieved using TRAVOS.

4.5.2.1 Calculating Trust

Using the information from Tables 4.2 and 4.3, a_1 can determine the number of successful interactions, m , and the number of unsuccessful interactions, n , for each agent it has interacted with. Feeding these into Equations 4.9 and 4.10, a_1 can obtain a parameter distribution which summarises each agent's likely behaviour in future interactions; for example, the shape parameters α and β , for a_2 , are calculated as follows:

Using Table 4.2: $m_{a_1,a_2} = 17$, $n_{a_1,a_2} = 5$.

Using Equations 4.9 & 4.10: $\alpha = 17 + 1 = 18$ and $\beta = 5 + 1 = 6$.

The hyperparameter for each agent is then used to estimate the probability that each agent will cooperate in any future interaction. In line with Section 4.1, we calculate this estimate as the expected value of the parameter distribution (Equation 4.5); for example, the estimate, ϑ_{a_1,a_2} , for a_2 is calculated as follows:

Using Equation 4.5: $\vartheta_{a_1,a_2} = \frac{\alpha}{\alpha+\beta} = \frac{18}{18+6} = 0.75$.

The above estimate gives a_1 an assessment of a_2 's likely behaviour based on direct interactions. However, as discussed in Section 4.4, a_1 may wish to determine if the accuracy of this estimate is sufficient to avoid the need to gather reputation. To do this, we calculate the posterior probability that the true value for θ_{a_1,a_2} lies within an acceptable

Agent	α	β	ϑ_{a_1, a_x}	γ_{a_1, a_x}
a_2	18	6	0.75	0.9806
a_3	3	16	0.1579	0.9798
a_4	19	6	0.76	0.9835
a_5	10	15	0.4	0.9657
a_6	4	1	0.8	0.8704
a_7	19	12	0.6129	0.9822

TABLE 4.4: Agent a_1 's calculated trust and associated confidence level for HTML content and phone call service provider agents.

margin of error around the estimate. We can calculate this using the parameter distribution as follows. First, we decide on an acceptable error margin, $\vartheta_{a_1 a_2} \pm \epsilon$, where ϵ is a suitable value, such as 0.2. Second, we integrate the parameter distribution over the area defined by the error margin. Finally, we determine some threshold for this probability, above which the estimate gives an acceptable level of accuracy; for example, we could define a threshold τ as 0.95. The proceeding example illustrates this calculation for a_1 's estimate for a_2 , using $\epsilon = 0.2$; we denote the resulting confidence value as γ_{a_1, a_2} :

$$\gamma_{a_1, a_2} = \frac{\int_{\vartheta_{a_1, a_2} - \epsilon}^{\vartheta_{a_1, a_2} + \epsilon} B^{\alpha-1} (1-B)^{\beta-1} dB}{\int_0^1 U^{\alpha-1} (1-U)^{\beta-1} dU} = \frac{\int_{0.55}^{0.95} B^{\alpha-1} (1-B)^{\beta-1} dB}{\int_0^1 U^{\alpha-1} (1-U)^{\beta-1} dU} = 0.98$$

The hyperparameters, estimate and associated confidence for each agent, a_2 to a_7 , which a_1 computes using TRAVOS, are shown in Table 4.4. From this, it is clear that the trust values for agents a_2 , a_3 and a_4 , all have a confidence above τ ($=0.95$). This means that a_1 does not need to consider the opinions of others for these three agents. Agent a_1 is able to decide that a_4 is the most trustworthy out of the three phone call service provider agents and chooses it to provide the phone call service for the VO.

4.5.2.2 Calculating Reputation

The process of selecting the most trustworthy HTML content service provider is not as straightforward. Agent a_1 has calculated that out of the possible HTML service providers, a_6 has the highest trust value. However, it has determined that the confidence it is willing to place in this value is 0.8704, which is below that of τ and means that a_1 has not yet interacted with a_6 enough times to calculate a sufficiently confident trust value. In this case, a_1 has to use the opinions from other agents that have interacted with a_6 , and form a reputation value for a_6 that it can compare to the trust values it has calculated for other HTML providers (a_5 and a_7).

Suppose that a_1 is aware of three agents that have interacted with a_6 , denoted by a_8 , a_9 and a_{10} , whose opinions about a_6 are $(15, 46)$, $(4, 1)$ and $(3, 0)$ respectively. Agent a_1 can then obtain hyperparameters based solely on the opinions provided as follows:

Opinions from providers: $a_8 = (15, 46)$, $a_9 = (4, 1)$ and $a_{10} = (3, 0)$

Using Equations 4.19 & 4.20: $M = 15 + 4 + 3 = 22$, $N = 46 + 1 + 0 = 47$

Using Equations 4.17 & 4.18: $\alpha = 22 + 1 = 23$, $\beta = 47 + 1 = 48$

Having obtained the shape parameters, a_1 can obtain an estimate for a_6 using Equation 4.5, as follows:

Using Equation 4.5: $\vartheta_{a_1, a_6} = \frac{\alpha}{\alpha + \beta} = \frac{23}{23 + 48} = 0.3239$

Now, a_1 is able to compare the trust in agents a_5 , a_6 and a_7 . Before calculating the trustworthiness of a_6 , agent a_1 considered a_6 to be the most trustworthy (see Table 4.4). Having calculated a new trust value for agent a_6 (which is lower than the first assessment), agent a_1 now regards a_7 as the most trustworthy. Therefore a_1 chooses a_7 as the service provider for the HTML content service.

4.5.2.3 Handling Inaccurate Opinions

The method a_1 uses to assess the trustworthiness of a_6 , as described in Section 4.5.2.2, is susceptible to errors caused by reputation providers giving inaccurate information. In our scenario, suppose a_8 provides the HTML content service too, and is in direct competition with a_6 . Agent a_1 is not aware of this fact, which makes a_1 unaware that a_8 may provide inaccurate information about a_6 to influence its decision on which HTML content provider agent to incorporate into the VO. If we examine the opinions provided by agents a_8 , a_9 and a_{10} , which are $(15, 46)$, $(4, 1)$ and $(3, 0)$ respectively, we can see that the opinion provided by a_8 does not correlate with the other two. Agents a_9 and a_{10} provide a positive opinion of a_6 , whereas agent a_8 provides a very negative opinion. Suppose that a_8 is providing an inaccurate account of its experiences with a_6 . We can use the mechanism discussed in Section 4.3 to allow a_1 to cope with this inaccurate information, and arrive at a better decision that is not influenced by self-interested reputation providing agents (such as a_8).

Before we show how TRAVOS can be used to handle such inaccurate information, we must assume the following. Agent a_1 obtained reputation information from a_8 , a_9 and a_{10} on several occasions, and each time a_1 recorded the opinion provided by a reputation provider and the actual observed outcome (from the interaction with an agent to which the opinion is applied). Each time an opinion is provided, the outcome observed is recorded by updating a frequency bin corresponding to the interval, Θ_r^C , which the received opinion belongs to. Agent a_1 maintains information of like opinions in bins as shown in Table 4.6. For example, if a_8 provides an opinion that is used to obtain a trust value of 0.254, then the actual observed outcome (successful or unsuccessful) is stored in the $0.2 < E[\theta_{a_{tr}, a_{te}} | \phi_r] \leq 0.4$ bin.

Agent	Weighting	Adjusted Values			
		μ	σ	α	β
a_8	0.0049	0.4988	0.2875	1.0095	1.0144
a_9	0.7802	0.6672	0.1881	3.5215	1.7567
a_{10}	0.7424	0.7227	0.1956	3.0629	1.1751

TABLE 4.5: Agent a_1 's adjusted values for opinions provided by a_8 , a_9 and a_{10} .

	[0, 0.2]		[0.2, 0.4]		[0.4, 0.6]		[0.6, 0.8]		[0.8, 1]		Total
	m	n	m	n	m	n	m	n	m	n	
a_8	2	0	11	4	0	0	0	0	2	3	22
a_9	0	2	1	3	0	0	22	10	6	4	48
a_{10}	1	3	0	2	0	0	18	8	5	3	40

TABLE 4.6: Observations made by a_1 given opinions from reputation sources. m represents that the interaction (to which the opinion applied) was successful, and likewise n means unsuccessful.

Using the information shown in Table 4.6, agent a_1 can calculate the weighting to be applied to the opinions from the three reputation sources by applying the technique described in Section 4.3. In so doing, agent a_1 uses the information from the bin that contains the opinion provided, and integrates the beta distribution between the limits defined by the bin's boundary. For example, a_8 's opinion falls under the $0.2 < E[\theta_{a_{tr}, a_{te}} | \phi_r] \leq 0.4$ bin. In this bin, agent a_1 has recorded that $m = 11$ and $n = 4$. These m and n values are used to obtain a beta distribution, $d(\theta_{a_{tr}, a_{te}} | \phi_o)$, which is then integrated between 0.2 and 0.4 to give a probability of accuracy $\rho_{a_1, a_8} = 0.0049$ for a_8 's opinion. Then, by using Equations 4.22 and 4.23, agent a_1 can calculate the adjusted mean and standard deviation of the opinion, which in turn gives the adjusted α and β parameters for that opinion. The results from these calculations are shown in Table 4.5.

Summing the adjusted values for α and β from Table 4.5, a_1 can obtain a more reliable value for the trustworthiness of a_6 . Using Equation 4.5, a_1 calculates an estimate $\vartheta_{a_1, a_6} = 0.7419$ for a_6 . This means that from the possible HTML content providers, a_1 now sees a_6 as the most trustworthy and selects it to be a partner in the VO. Unlike a_1 's decision in Section 4.5.2.2 (when a_7 was chosen as the VO partner), here we have shown how a reputation provider cannot influence the decision made by a_1 by providing inaccurate information.

4.6 Empirical Study

In this section, we demonstrate the advantages that TRAVOS offers to the state of the art, through empirical evaluation. We divide our discussion into three parts. First, Section 4.6.1 describes the simulation environment and overall methodology used to perform our experiments. Second, Section 4.6.2 compares the reputation component of

TRAVOS to the Beta Reputation System (BRS) (see Sections 2.2.2 & 2.3.2 for more detail). We have chosen this model as a benchmark because it shares the same basic representation of trust as TRAVOS. Any difference in performance can therefore be attributed to the novel properties of TRAVOS rather than those it shares with BRS. Finally, Section 4.6.3 investigates the component performance of TRAVOS; that is, how TRAVOS performs when a truster uses both its direct experience of a trustee and reputation, and when it uses either source of evidence in isolation. This allows us to show how TRAVOS behaves when different types of information are available, and that using both types of information is in general better than using one or the other independently.

4.6.1 Experiment Methodology

Evaluation of TRAVOS took place using a simulated marketplace environment, consisting of three distinct sets of agents: provider agents $\mathcal{P} \subset \mathcal{A}$, consumer agents $\mathcal{C} \subset \mathcal{A}$, and reputation source agents $\mathcal{S} \subset \mathcal{A}$. For our purposes, the role of any $c \in \mathcal{C}$ is to evaluate $\vartheta_{c,p}$ for all $p \in \mathcal{P}$. The behaviour of each provider and reputation source agent was set before each experiment. Specifically, the behaviour of a provider $p_1 \in \mathcal{P}$ is determined by the parameter θ_{c,p_1} as described in Section 3.3. Here, reputation sources are divided into three types that define their behaviour: *accurate* sources report the number of successful and unsuccessful interactions they have had with a given consumer without modification; *noisy* sources add Gaussian noise to the beta distribution determined from their interaction history, rounding the resulting expected value if necessary to ensure that it remains in the interval $[0, 1]$; and *lying* sources attempt to maximally mislead the consumer by setting the expected value E^r to $1 - E^r$.

Against this background, all experiments consisted of a series of episodes in which a consumer was asked to assess its trust in all providers \mathcal{P} . Based on these assessments, we calculate the consumer's mean estimation error for the episode (Equation 4.30). This gives us a measure of the consumer's performance on assessing the provider population as a whole. The value of this metric will vary depending on the distribution of values of $\theta_{c,p}$ over the provider population. For simplicity, all the results described in the next sections have been acquired for a population of 101 providers, $\{p_1, \dots, p_n\}$, with values of $\theta_{c,p}$ chosen uniformly between 0 and 1 at intervals of 0.01.

$$avg_estimate_err = \frac{1}{n} \sum_{i=1}^n abs(\vartheta_{c,p_i} - \theta_{c,p_i}) \quad (4.30)$$

In each episode, the consumer may draw upon both the opinions of reputation sources in \mathcal{S} and its own interaction history with both the providers and reputation sources. However, to ensure that the results of each episode are independent, the interaction history between all agents is cleared before every episode, and re-populated according to

experiment	no. lying	no. noisy	no. accurate
1	0	0	20
2	0	10	10
3	0	20	0
4	10	0	10
5	20	0	0

TABLE 4.7: Reputation source populations.

set parameters. All the results that we discuss have been tested for statistical significance using Analysis of Variance techniques and Scheffé tests (Cohen, 1995).

4.6.2 TRAVOS Against the Beta Reputation System

Like TRAVOS, BRS uses the beta family of probability functions to calculate the posterior probability of an agent a_{te} 's behaviour holding a certain value, given past interactions with a_{te} . However, the models differ significantly in their approach to handling inaccurate reputation. TRAVOS assesses each reputation source individually, based on the perceived accuracy of past opinions. In contrast, BRS assumes that the majority of reputation sources provide an accurate opinion, and it ignores any opinions that deviate significantly from the average. Since BRS does not differentiate between reputation and direct observations, we have focused our evaluation on scenarios where consumers have no personal experience, and must therefore rely on reputation only.

To show variation in performance depending on reputation source behaviour, we ran experiments with populations containing accurate and lying reputation sources, and populations containing accurate and noisy sources. In each case, we kept the total number of sources equal to 20, but ran separate experiments in which the percentage of accurate sources was set to 0%, 50% and 100% (see Table 4.7). Figure 4.6 shows the mean estimation error of TRAVOS and BRS with these different reputation source populations averaged over 50 independent episodes in each experiment. To provide a benchmark, the figure also shows the mean estimation error of a consumer, $c_{0.5}$, which keeps $\vartheta_{c_{0.5},p} = 0.5$ for all $p \in \mathcal{P}$. Results are plotted against the number of previous interactions between the consumer and each reputation source.

As can be seen, in populations containing lying agents, the mean estimation error of TRAVOS is consistently equal to or less than that of BRS. Moreover, estimation errors decrease significantly for TRAVOS as the number of consumer to reputation source interactions increases. In contrast, BRS's performance remains constant, since it does not learn from past experience. Both models perform consistently better than $c_{0.5}$ in populations containing 50% or 0% liars. However, in populations containing only lying sources, both models are sufficiently misled to perform worse than $c_{0.5}$, but TRAVOS suffers less from this effect than BRS. Specifically, when the number of past consumer

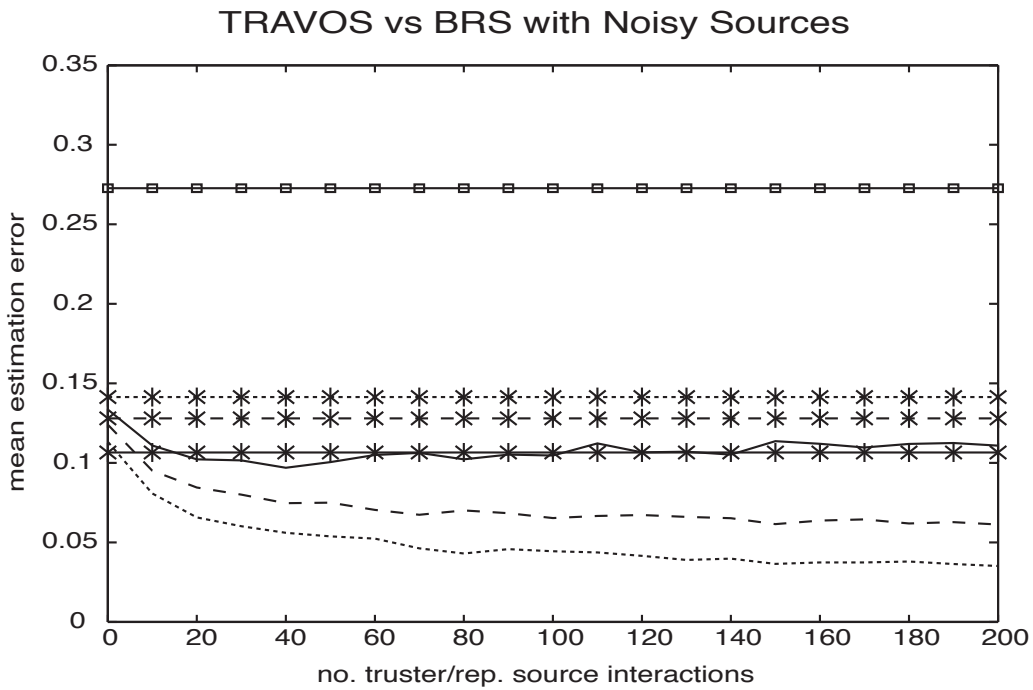
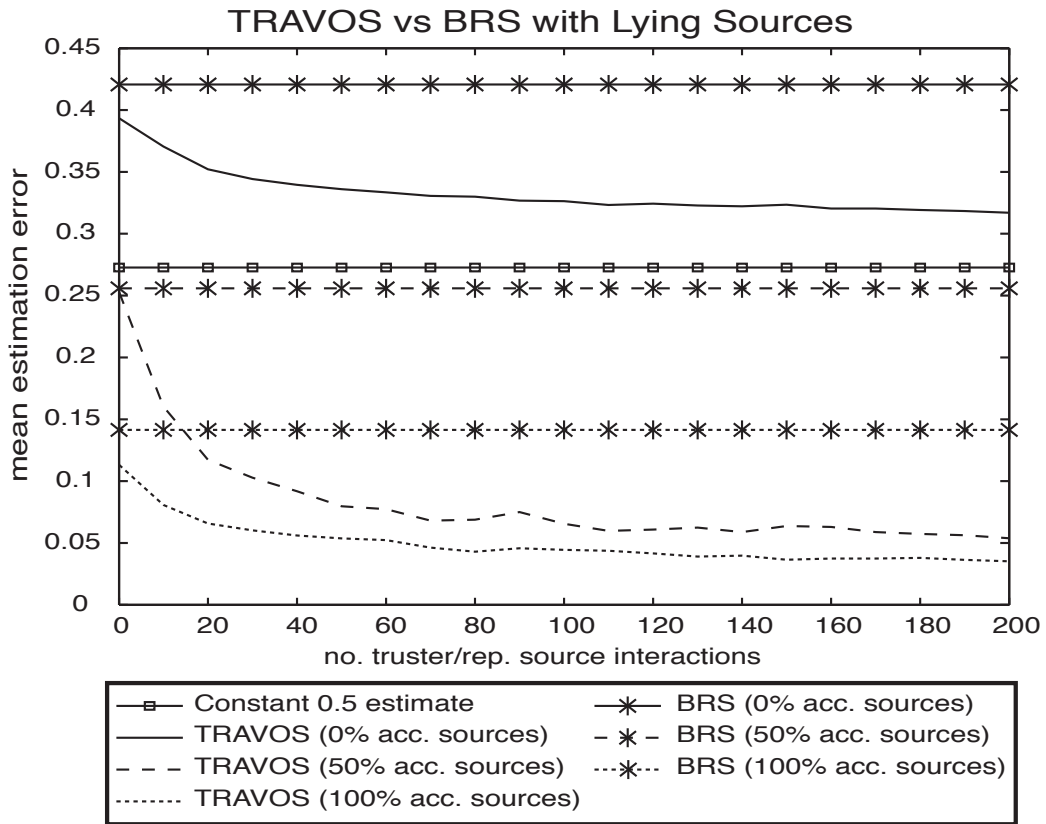


FIGURE 4.6: TRAVOS reputation system versus BRS.

to reputation interactions is low, TRAVOS benefits from its initially conservative belief in reputation source opinions. The benefit is enhanced further as the consumer becomes more skeptical with experience.

Similar results can be seen in populations containing noisy sources. In general, performance is better because noisy source opinions are not as misleading as lying source opinions on average. TRAVOS still outperforms BRS in most cases, except when the population contains only noisy sources. In this case, BRS has a small but statistically significant advantage when the number of consumer to reputation source interactions are less than 10.

4.6.3 TRAVOS Component Performance

To evaluate the overall performance of TRAVOS, we compared three versions of the system that used the following information respectively: direct interactions between the consumer and providers; direct provider experience and reputation; and reputation information only. In these experiments, we varied the number of interactions between the consumers and providers, and kept the number of consumer to reputation source interactions constant at 10. We used the same reputation source populations as described in Section 4.6.2. The mean estimation errors for a subset of these experiments are shown in Figure 4.7. Using only direct consumer to provider experience, the mean estimation error decreases as the number of consumer to provider interactions increases. As would be expected, using both information sources when the number of consumer to provider interactions is low results in similar performance to using reputation information only. However, in some cases, the combined model may provide marginally worse performance than using reputation only.⁷ This can be attributed to the fact that TRAVOS always puts more faith in direct experience than reputation.

With a population of 50% lying reputation sources, the combined model is misled enough to temporarily increase its error rate above that of the direct only model. This is a symptom of the relatively small number of consumer to reputation source interactions (10), which is insufficient for the consumer to completely discount all the reputation information as unreliable. The effect disappears when the number of such interactions is increased to 20. However, these results are not illustrated graphically here.

4.7 Summary

In this chapter, we introduce the TRAVOS trust model, which instantiates the framework introduced in the previous chapter, for cases in which a trustee's behaviour can

⁷This effect was not considered significant under a Scheffé test, but was considered significant by Least Significant Difference Testing. The latter technique is, in general, less conservative at concluding that a difference between groups does exist.

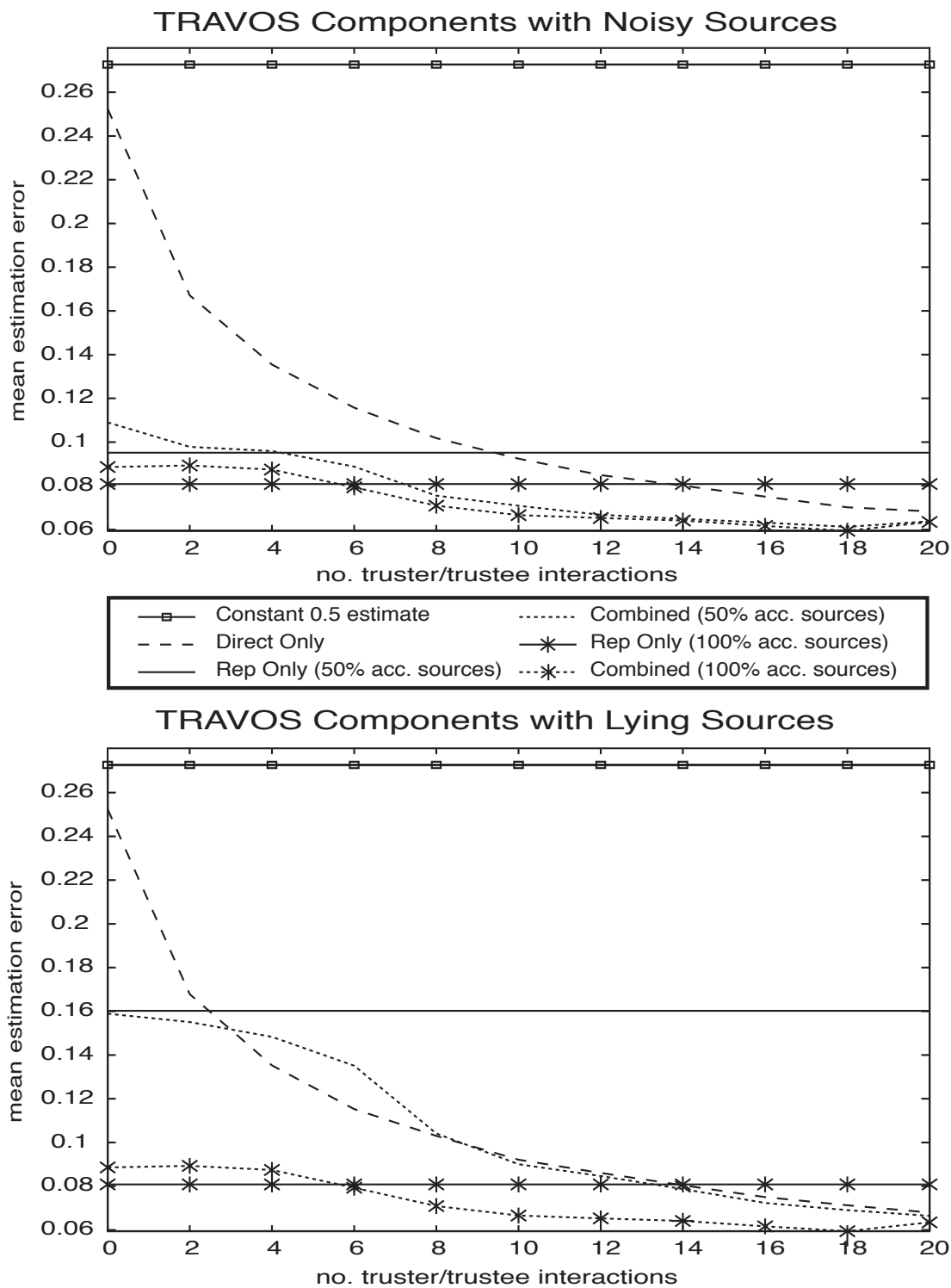


FIGURE 4.7: TRAVOS component performance.

be described as a binary event. That is, TRAVOS deals with situations in which a truster only cares about whether or not a trustee cooperates by fulfilling its obligations, or defects by breaking those obligations.

In particular, we showed how, by using TRAVOS, a truster could decide if and when to interact with a trustee based on its direct observations and reputation. Significantly, as part of this solution, we introduced a reputation filtering mechanism that allows a truster to account for the reliability of its reputation sources in making its assessment. To achieve this, the truster first estimates the probability with which a source's prediction is within a margin of error around a trustee's true behaviour. This is then used as part of a heuristic method, to reduce the impact of unreliable sources on the truster's decisions.

Following on from this, we then demonstrated how TRAVOS works in practice in two ways. First, we showed how it can be used as part of the CONOISE-G system to guide automatic formation and management of agent-based virtual organisations. This was followed by a walk-through scenario, showing how the model parameters respond to observations made in the environment. Second, we demonstrated the performance of TRAVOS through empirical analysis. In particular, this showed how it out-performs the most similar model in the literature, under most of the conditions tested.

Chapter 5

TRAVOS-C: A Trust Model for Continuous Action Spaces

In the previous chapter, we introduced the TRAVOS system, and showed how it could be used to fulfill our aim of facilitating decision making in a multi-agent system. As discussed, however, TRAVOS relies on a heuristic method for filtering out inaccurate opinions, and is only applicable when trustees have binary action spaces.

In this chapter, we improve upon this situation by introducing a refinement, TRAVOS-C¹, which not only addresses continuous action spaces, but has four other key advantages over its predecessor. First, reputation source accuracy is no longer assessed by a heuristic, but follows directly from Bayesian theory along with the model's assumptions. This gives the model a more solid theoretical foundation that is optimal under the model's assumptions.

Second, in addition to assessing the accuracy of a reputation source based on its past performance, the Bayesian model can account for observed correlations between opinions from different sources. This means that, if a group of agents generally provide similar opinions, then evidence for the reliability of one of the group can count against the group as a whole.

Third, if a reputation source always provides biased opinions, TRAVOS-C can still make use of the source, provided that a correlation with trustee behaviour still exists. For example, if each time Alice gives Bob an apple, Bob says that she gave him an orange, then knowing that Bob reports receiving an orange is evidence that he actually received an apple. Under TRAVOS, such reports would be ignored, even if there was an obvious correlation.

Finally, the model can further improve the accuracy of its estimates, by observing correlations between the behaviour of groups of agents. For example, if it is observed that

¹The C in the name refers to the applicability of the new model to continuous action spaces.

agents belonging to a particular organisation generally provide a certain quality of service, the truster can use this information to better estimate a trustee's behaviour when both direct observations and reputation are in short supply.

In the following sections, we elaborate on these claims and detail the theoretical basis for TRAVOS-C. We begin in Section 5.1 by defining the basic model and outlining how it can be used to assess trustee behaviour based on direct experience and reputation. Building on this, Section 5.2 shows how the model can be extended to account for correlations in group behaviour, and Section 5.3 instantiates the model for continuous outcome spaces. Finally, Sections 5.4 and 5.5 detail how the model can be applied in practice, Section 5.6 gives an empirical evaluation of the model, and Section 5.7 summarises.

5.1 The TRAVOS-C Model

As in previous chapters, our main goal is to estimate the future behaviour distribution for a trustee, given direct and third party observations of agent behaviour. In the case of third party information, different reputation sources may be more reliable than others, so a truster must decide how much influence each reputation source should have on trustee assessment. While making this decision, the important thing to consider is the relative predictive value of opinions from different sources. However, provided we can assess their relative worth, the reasons why one reputation source provides better opinions than another are only of secondary importance. For instance, if opinions from a particular agent generally have low correlation with trustee behaviour, this may be due to the reputation source purposely misleading the truster, or because the reputation source has an inaccurate world view. In both cases, the effect is the same: the reputation source's opinions provide little information about trustee behaviour, and so should not have a significant influence on trustee assessment.

This gives the intuition behind TRAVOS-C, in which *any* source of inaccuracy in third party opinions is modelled as independent random noise added to each reputation source's observations of a trustee, before an opinion is conveyed to the truster. Thus, an unreliable reputation source is modelled as having significant noise associated with its opinions, which, as result, provide little information to the truster. In contrast, a reliable reputation source is modelled with little or no added noise. Thus, its opinions will be judged to have value similar to direct experience. In either case, Bayesian inference is applied to determine the probable amount of noise contained in a reputation source's opinion. This replaces the filtering mechanism used in TRAVOS (Section 4.3), whereby the weight of an opinion is reduced heuristically, if evidence suggests that it deviates significantly from trustee behaviour.

The formal aspects of this approach are illustrated by the Bayesian network in Figure 5.1, in which the dashed ovals indicate which variables are assumed to be visible

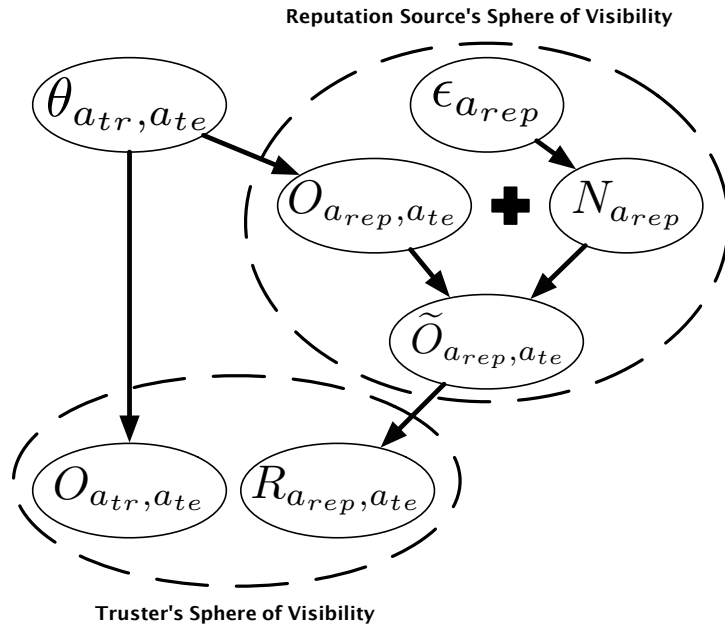


FIGURE 5.1: The TRAVOS-C model.

to the truster and each of its reputation sources respectively. As before, the behaviour distribution of a trustee toward a truster is determined by a parameter vector denoted $\theta_{a_{tr}, a_{te}}$; the interaction history between the trustee and truster is denoted $O_{a_{tr}, a_{te}}^{0:t'}$; and the interaction history between the trustee and reputation source is denoted $O_{a_{rep}, a_{te}}^{0:t'}$. As indicated in the figure, the truster can only observe its own direct observations and the reported opinions of reputation sources. If any third party opinion is to be useful, then it should depend on the observation history between trustee and reputation source. However, TRAVOS-C introduces two modifications to the basic framework. First, we assume that the behaviour distribution of the trustee toward the reputation source is *always* identical to its behaviour distribution toward the truster. This implies that $O_{a_{rep}, a_{te}}^{0:t'}$ is a private set of observations drawn from a_{te} 's behaviour distribution toward a_{tr} that, in a perfect world, will always be returned unmolested to the truster. Second, we introduce a new random variable $N_{a_{rep}}$, which has a probability distribution specified by a parameter vector $\epsilon_{a_{rep}}$, and a domain denoted \mathcal{N}^C . In the model, $N_{a_{rep}}$ plays the role of noise added to each $O_{a_{rep}, a_{te}}$ to form noisy observations, denoted $\tilde{O}_{a_{rep}, a_{te}}$. To ensure that both noisy and direct observations can be treated in the same way, we define \mathcal{N}^C such that:

$$\forall N_{a_{rep}} \in \mathcal{N}^C, \forall O_{a_{rep}, a_{te}} \in \mathcal{O}^C, \tilde{O}_{a_{rep}, a_{te}} = (N_{a_{rep}} + O_{a_{rep}, a_{te}}) \in \mathcal{O}^C \quad (5.1)$$

With this in mind, we define $\tilde{O}_{a_{rep}, a_{te}}^{0:t'}$ as the set of all noisy observations obtained by a_{rep} up to time t' , and it is on this set that we make $R_{a_{rep}, a_{te}}$ a statistic.

Using this model, practical Bayesian inference can be performed by considering the posterior distribution of $\theta_{a_{tr}, a_{te}}$ given $O_{a_{tr}, a_{te}}^{0:t'}$ and $R_{a_{rep}, a_{te}}$ for a number of reputation

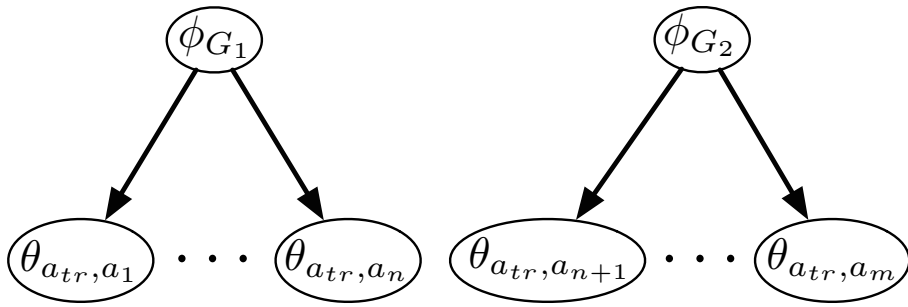


FIGURE 5.2: Bayesian network for inferring group priors.

sources a_{rep} . Although details of how this can be done are deferred to the following sections, there are essentially two mechanisms by which this model can provide effective predictions of behaviour. The first and simplest of these is direct interaction, whereby the truster builds up a picture of the trustee's behaviour by interacting with it repeatedly over time.

The second is similar to how reputation is used in TRAVOS, in which the reliability of a reputation source is assessed by receiving opinions about different trustees over time. This works by identifying plausible levels of noise that explain the amount of correlation observed between opinions and later observed direct observations; the direct observations provide reliable information about trustee behaviour, and any discrepancy between this and reported opinions points to high levels of noise. In turn, if higher levels of noise seem plausible, then the reputation source in question will have little effect on the posterior distributions of trustee behaviour.

5.2 Learning Group Behaviour

So far we have discussed models for assessing trustee behaviour when either the truster or some third party has experience of the trustee. However, there are situations where this may not be the case, for example when a new service provider enters a system for the first time with no previous provision history. One way to deal with such cases is to draw comparisons between agents with certain attributes in common. For instance, if we find that agents that belong to a particular organisation generally behave in a particular way then, *a priori*, we could assume that an unknown agent from that group will behave similarly to its peers.

To take advantage of such correlations, TRAVOS-C includes an additional component that can be invoked when correlations between groups of agents are likely to exist. For instance, using the multimedia scenario from Section 4.5.2, suppose that a number of agents provide multimedia content using the same type of streaming technology. Presumably, the quality of this technology will be a major factor in the capabilities of

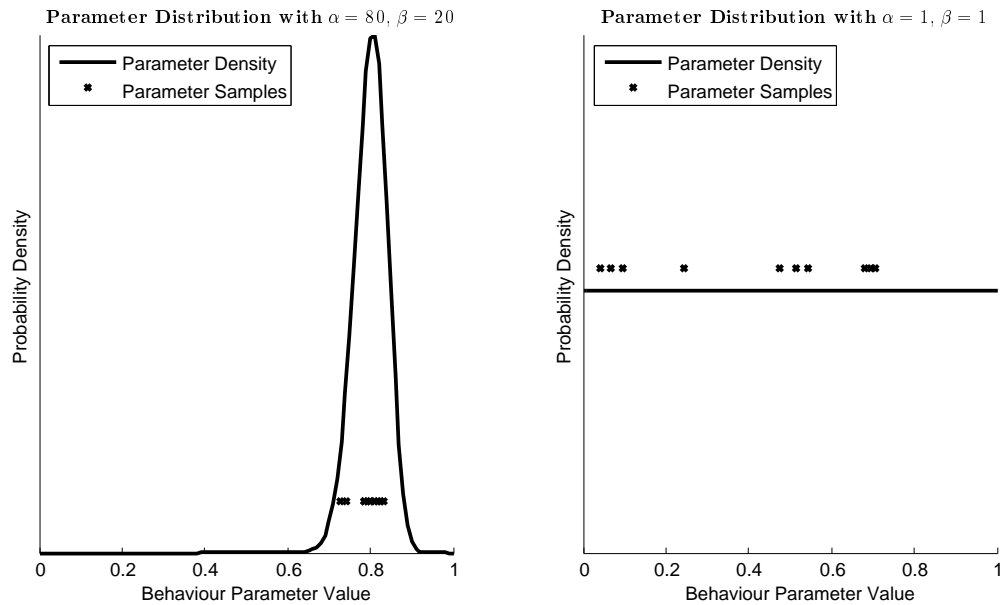


FIGURE 5.3: Examples of group parameter distributions with behaviour samples.

these agents to provide their service. Therefore, it is not unreasonable to speculate that the ability of one such agent to deliver multimedia will be similar to that of another. If such correlations do exist, then we can use the behaviour of known agents to predict the behaviour of other agents that use the same technology.

The idea, then, is to divide agents into groups containing members that we suspect will have similar behaviours; for example, this could be based on organisation membership, or geographical location. Then, we associate with each group a parameter distribution, which characterises the distribution of behaviours belonging to agents in that group. We can use this to determine things such as the average behaviour of agents belonging to a particular group and how similar agent behaviours are within a group, and to predict the behaviour of an anonymous agent within a group. Put formally, we separate the set of agents \mathcal{A} into distinct subsets, G_1, \dots, G_g , such that $\bigcup_{i=1}^g G_i = \mathcal{A}$. Then we associate a hyperparameter vector, ϕ_{G_i} , with each group G_i , which determines the parameter distribution of trustee behaviours belonging to G_i .

The way this model works is demonstrated in Figures 5.2 and 5.3. Specifically, Figure 5.2 shows a Bayesian network for two distinct sets of trustees, $G_1 = \{a_1, \dots, a_n\}$ and $G_2 = \{a_{n+1}, \dots, a_m\}$. Here, the parameter vectors for all the members of each group are shown to be dependent on their respective hyperparameter vector. This is exemplified in Figure 5.3, in which the parameter distributions for each group are depicted along with trustee behaviour distributions sampled from them. Here, we assume that trustees have binary action spaces for illustration purposes. In this case, each parameter vector, $\theta_{a_{tr}, a_{te}}$, specifies the probability of the trustee supplying a successful service, as is assumed in TRAVOS. The group parameter distributions become beta distributions, with α and β parameters specified by each hyperparameter ϕ_{G_i} . For G_1 , the hyperparameter

ϕ_{G_1} specifies $\alpha = 80$, $\beta = 20$. This is an example of a highly informative group, in which member behaviours do not deviate significantly from a mean of 0.8. On the other hand, ϕ_{G_2} gives $\alpha = 1$ and $\beta = 1$, which gives a non-informative uniform distribution. In this case, knowing that an agent belongs to G_2 tells us relatively little, other than that a member of G_2 is equally likely to assume any behaviour distribution.

The application of this model is highly flexible. For instance, we could potentially assign agents to groups based on organisation membership, social relationships, or the length of time an agent has been in the system. Obviously, it is best to pick predictive attributes based on some investigation of the target application. However, choosing relatively uninformative attributes will not reduce the effectiveness of the system because the model will automatically account for how informative a group is. Even if no predictive attributes can be identified, the group model may still provide some value by assigning all agents in an environment to a single group. In this case, TRAVOS-C could learn the most appropriate prior to use for the environment as a whole. In this way, if most agents in a system behave in a certain way, TRAVOS-C will use this information to give a head start in behaviour prediction, and so will generally need fewer interactions with a trustee to accurately predict its behaviour.

In addition, the method by which group parameter distributions are determined is highly flexible. For example, we could statically assign group parameter distributions, or more significantly, apply standard Bayesian techniques to learn appropriate hyperparameters dynamically, by observing correlations in trustee behaviour within groups. In fact, committing to a fixed hyperparameter vector per group is not even strictly necessary. Instead, we can consider the joint distribution of both the hyperparameters and the behaviour parameters alike, and marginalise over all possible hyperparameters. This accounts for any uncertainty surrounding the hyperparameters themselves.

For example, if a truster interacts ten times with two distinct members of a group and finds that, over those interactions, the trustees appear to behave in a similar way, this may suggest a highly informative group parameter distribution. On the other hand, perhaps ten interactions are not enough to pin down the individual behaviour distributions, and perhaps two trustees are not a representative sample of group members. However, by marginalising over the hyperparameters, we need not worry about such concerns, because the Bayesian machinery will always give the most appropriate distribution for trustee behaviour. That is, the marginal distribution will account for both the uncertainty inherent in the trustee's behaviour, and the uncertainty due to lack of evidence. By using this distribution, we can then make choices using Decision Theory, to account appropriately for the risk due to both sources of uncertainty. In the following sections we shall see how this is achieved for a particular representation of trustee behaviour.

5.3 Instantiating TRAVOS-C for Continuous Action Spaces

At this level of detail, the model could potentially be applied to any domain of trustee behaviour², including continuous and binary action spaces. However, to fully assess the merits of this approach requires a fully instantiated working model, and this necessitates the choice of some set of assumptions to allow learning to take place. For this reason, and since we wish to demonstrate applicability to continuous action spaces, we assume that both trustee behaviour and reputation noise are normally distributed.

This is not an unreasonable assumption, because there are many instances where real life phenomena have distributions that are approximately normal. For example, according to the central limit theorem, many random variables that can be seen as the sum of several other random variables will have an approximately normal distribution. Despite this choice, however, we stress that the general methods discussed here could be applied to any reasonable choice of distribution.

With this mind, we must instantiate three aspects of the model: (1) the domains of each model parameter, (2) the reputation function for forming opinions, and (3) the prior distributions for each parameter to enable Bayesian inference and group behaviour analysis. We consider each of these in turn in the subsections that follow.

5.3.1 The Parameter Domains

Starting with the model parameters, the support³ of a normal distribution is the entire real line, so we let $\mathcal{O}^C = \mathcal{N}^C = \mathbb{R}$. This clearly satisfies Equation 5.1, since \mathbb{R} is closed under addition. Furthermore, as normal distributions are fully defined by their mean, μ , and variance, σ^2 , we let $\theta_{a_{tr}, a_{te}} = \langle \mu_{a_{tr}, a_{te}}, \sigma_{a_{tr}, a_{te}}^2 \rangle$ and $\epsilon_{a_{rep}} = \langle \mu_{\epsilon_{a_{rep}}}, \sigma_{\epsilon_{a_{rep}}}^2 \rangle$, where in each case, the subscripts identify which distribution the parameter belongs to. Accordingly, the domain of $\theta_{a_{tr}, a_{te}}$ becomes $\Theta^C = \mathbb{R} \times \mathbb{R}^+$, since μ is real and σ^2 is positive real.

5.3.2 The Reputation Function

From Chapter 3 we know that a reputation function should ideally be concise, but still capture all the information conveyed by the noisy observations on which an opinion is supposedly based. In more formal terms, we want to ensure that the model parameters, comprising $\theta_{a_{tr}, a_{te}}$ and each $\epsilon_{a_{rep}}$, are conditionally independent of all noisy observations,

²In some cases, this may require introducing a dependency between $N_{a_{rep}}$ and $O_{a_{rep}, a_{te}}$ so that Equation 5.1 can be satisfied. However, we do not concern ourselves with such cases here.

³The support of a probability distribution is the smallest closed set whose complement has probability zero.

given the third party opinions. Symbolically, this means that, for reputation sources a_1, \dots, a_n :

$$\begin{aligned} p(\theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(1)}}, \dots, \epsilon_{a_{rep(n)}} | R_{a_1, a_{te}}, \dots, R_{a_n, a_{te}}, \tilde{O}_{a_1, a_{te}}^{0:t'}, \dots, \tilde{O}_{a_n, a_{te}}^{0:t'}) = \\ p(\theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(1)}}, \dots, \epsilon_{a_{rep(n)}} | R_{a_1, a_{te}}, \dots, R_{a_n, a_{te}}) \end{aligned} \quad (5.2)$$

which, assuming an opinion is derived only from the noisy observations, is equivalent to:

$$\begin{aligned} p(\theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(1)}}, \dots, \epsilon_{a_{rep(n)}} | \tilde{O}_{a_1, a_{te}}^{0:t'}, \dots, \tilde{O}_{a_n, a_{te}}^{0:t'}) = \\ p(\theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(1)}}, \dots, \epsilon_{a_{rep(n)}} | R_{a_1, a_{te}}, \dots, R_{a_n, a_{te}}) \end{aligned} \quad (5.3)$$

From Bayes rule, and assuming that noise added by different sources is independent⁴, we know that

$$\begin{aligned} p(\theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(1)}}, \dots, \epsilon_{a_{rep(n)}} | \tilde{O}_{a_1, a_{te}}^{0:t'}, \dots, \tilde{O}_{a_n, a_{te}}^{0:t'}) \propto \\ p(\theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(1)}}, \dots, \epsilon_{a_{rep(n)}}) \prod_{i=1}^n p(\tilde{O}_{a_i, a_{te}}^{0:t'} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(i)}}) \end{aligned} \quad (5.4)$$

and that

$$\begin{aligned} p(\theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(1)}}, \dots, \epsilon_{a_{rep(n)}} | R_{a_1, a_{te}}, \dots, R_{a_n, a_{te}}) \propto \\ p(\theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(1)}}, \dots, \epsilon_{a_{rep(n)}}) \prod_{i=1}^n p(R_{a_i, a_{te}} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(i)}}) \end{aligned} \quad (5.5)$$

Since these equations only differ by their likelihood terms, $p(\tilde{O}_{a_i, a_{te}}^{0:t'} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(i)}})$ and $p(R_{a_i, a_{te}} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(i)}})$, we can prove that Equation 5.3 holds if we can prove that:

$$\forall i \in [1, n], p(\tilde{O}_{a_i, a_{te}}^{0:t'} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(i)}}) = p(R_{a_i, a_{te}} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep(i)}}) \quad (5.6)$$

To achieve this, we suggest defining the reputation function such that $R_{a_{rep}, a_{te}}$ comprises three values: (1) the number of observations in $\tilde{O}_{a_{rep}, a_{te}}^{0:t'}$ (denoted n), (2) the sample mean of $\tilde{O}_{a_{rep}, a_{te}}^{0:t'}$ (denoted m), and (3) its sample variance (denoted v). Thus, for each a_{rep} and a_{te} , we define $R_{a_{rep}, a_{te}}$ as the vector $\langle n, m, v \rangle$, where:

$$n = |\tilde{O}_{a_{rep}, a_{te}}^{0:t'}| \quad (5.7)$$

$$m = \frac{1}{n} \sum_{i=1}^n \tilde{O}_{a_{rep}, a_{te}}^i \quad (5.8)$$

$$v = \frac{1}{n} \sum_{i=1}^n (\tilde{O}_{a_{rep}, a_{te}}^i - m)^2 \quad (5.9)$$

$$(5.10)$$

⁴This can safely be assumed, provided reputation sources do not share information during opinion formation.

To prove that this satisfies our conditions, we first note that each noisy observation, $\tilde{O}_{a_{rep}, a_{te}}$, is a sum of two normally distributed random variables, so is itself normally distributed, with mean $\mu_{a_{tr}, a_{te}} + \mu_{\epsilon_{a_{rep}}}$ and variance $\sigma_{a_{tr}, a_{te}}^2 + \sigma_{\epsilon_{a_{rep}}}^2$. This means that $\tilde{O}_{a_{rep}, a_{te}}^{0:t'}$ is a set of independent normally distributed random variables with shared mean and variance, and we want to show that its likelihood, $p(\tilde{O}_{a_{rep}, a_{te}}^{0:t'} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep}})$, depends only on the size of the set, its mean and its variance. In the following sections, we shall encounter this problem on a number of occasions, as well as more general cases where samples have a shared mean, but possibly different variances. Therefore, we state the more general proof in Theorem 5.1 for future reference, and address the specific case of shared variance in Corollary 5.2. By applying these general results to the model, we can state that

$$\begin{aligned} p(\tilde{O}_{a_{rep}, a_{te}}^{0:t'} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep}}) \\ = p(R_{a_{rep}, a_{te}} | \theta_{a_{tr}, a_{te}}, \epsilon_{a_{rep}}) \end{aligned} \quad (5.11)$$

$$= \left[2\pi(\sigma_{a_{tr}, a_{te}}^2 + \sigma_{\epsilon_{a_{rep}}}^2) \right]^{-n/2} \cdot \exp \left[-\frac{n \left((m - \mu_{a_{tr}, a_{te}} - \mu_{\epsilon_{a_{rep}}}) + v \right)^2}{2(\sigma_{a_{tr}, a_{te}}^2 + \sigma_{\epsilon_{a_{rep}}}^2)} \right] \quad (5.12)$$

Theorem 5.1 (Gaussian Likelihood). *Suppose that $X = \{x_1, \dots, x_n\}$ is a set of n independent samples drawn from gaussian distributions with the same mean, μ , but possibly different variances, $\{\sigma_1^2, \dots, \sigma_n^2\}$. The full data likelihood function is then:*

$$\mathcal{L}(X) = (2\pi)^{-n/2} \cdot \exp \left[-\frac{\tau_{tot}}{2} \left((\bar{x} - \mu)^2 + s^2 \right) \right] \prod_{i=1}^n \frac{1}{\sqrt{\sigma_i^2}} \quad (5.13)$$

where

$$\tau_{tot} = \sum_{i=1}^n \frac{1}{\sigma_i^2} \quad (5.14)$$

$$\bar{x} = \frac{1}{\sum_{j=1}^n \frac{1}{\sigma_j^2}} \sum_{i=1}^n \frac{x_i}{\sigma_i^2} \quad (5.15)$$

$$s^2 = \left(\frac{1}{\sum_{j=1}^n \frac{1}{\sigma_j^2}} \sum_{i=1}^n \frac{x_i^2}{\sigma_i^2} \right) - \bar{x}^2 \quad (5.16)$$

Proof:

It follows directly from the definition of the gaussian distribution that:

$$\mathcal{L}(X) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left[-\frac{(x_i - \mu)^2}{2\sigma_i^2} \right] \quad (5.17)$$

$$= (2\pi)^{-n/2} \cdot \exp \left[-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma_i^2} \right] \prod_{i=1}^n \frac{1}{\sqrt{\sigma_i^2}} \quad (5.18)$$

Now, if we expand the summation in Equation 5.18 we get:

$$\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma_i^2} = \sum_{i=1}^n \frac{x_i^2 - 2\mu x_i + \mu^2}{\sigma_i^2} \quad (5.19)$$

$$= \left(\sum_{i=1}^n \frac{x_i^2}{\sigma_i^2} \right) - 2\mu \left(\sum_{i=1}^n \frac{x_i}{\sigma_i^2} \right) + \mu^2 \left(\sum_{i=1}^n \frac{1}{\sigma_i^2} \right) \quad (5.20)$$

$$= \tau_{tot} \left[\left(\frac{1}{\tau_{tot}} \sum_{i=1}^n \frac{x_i^2}{\sigma_i^2} \right) - 2\mu\bar{x} + \mu^2 \right] \quad (5.21)$$

and by completing the square we have:

$$\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma_i^2} = \tau_{tot} \left[\left(\frac{1}{\tau_{tot}} \sum_{i=1}^n \frac{x_i^2}{\sigma_i^2} \right) - \bar{x}^2 + (\bar{x} - \mu)^2 \right] \quad (5.22)$$

$$= \tau_{tot} [s^2 + (\bar{x} - \mu)^2] \quad (5.23)$$

Finally, by substituting Equation 5.23 into Equation 5.18 we obtain Equation 5.13, thus proving the theorem.

Corollary 5.2 (Shared Variance Gaussian Likelihood). *If $X = x_1, \dots, x_n$ is a set of n independent samples drawn from a gaussian distribution with fixed mean, μ and variance, σ^2 , then it follows from Theorem 5.1 that the full data likelihood is:*

$$\mathcal{L}(X) = (2\pi\sigma^2)^{-n/2} \cdot \exp \left[-\frac{\tau_{tot}}{2} ((\bar{x} - \mu)^2 + s^2) \right] \quad (5.24)$$

where

$$\tau_{tot} = \frac{n}{\sigma^2} \quad (5.25)$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (5.26)$$

$$s^2 = \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right) - \bar{x}^2 \quad (5.27)$$

5.3.3 Parameter Distributions for Group Behaviour

The final part of TRAVOS-C that we must instantiate is the group behaviour model. That is, for a given trustee $a_{te} \in G_i$, we must specify the distribution of $\theta_{a_{tr}, a_{te}}$ given ϕ_{G_i} , which we can then use to derive the full posterior distribution of $\theta_{a_{tr}, a_{te}}$ given all available evidence. Since we postulate that trustee behaviour is normally distributed, it seems reasonable that we should define ϕ_{G_i} in line with the conjugate family for the Gaussian distribution. This generally keeps the form of equations as simple as possible, particularly when we ignore reputation and use only direct observations, in which case the posterior distribution of $\theta_{a_{tr}, a_{te}}$ inherits the same form as the prior (see Chapter 3).

Consequently, we model the group parameter distribution as a normal inverse-gamma distribution (Denison et al., 2002), which is well known to be the conjugate family for Gaussian distributions. In general, this has four hyperparameters, denoted m , v , α and β , which specify the p.d.f. of Gaussian parameters μ and σ^2 as follows:

$$p(\mu, \sigma^2 | m, v, \alpha, \beta) = p(\mu | \sigma^2, m, v) p(\sigma^2 | \alpha, \beta) \quad (5.28)$$

$$p(\mu | \sigma^2, m, v) = \frac{1}{\sqrt{2\pi v \sigma^2}} \exp \left[-\frac{(\mu - m)^2}{2v \sigma^2} \right] \quad (5.29)$$

$$p(\sigma^2 | \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} \exp \left[-\frac{\beta}{\sigma^2} \right] \quad (5.30)$$

Here, the distribution is defined in terms of the conditional probability of μ given σ^2 , and the marginal distribution of σ^2 . Specifically, μ has a normal distribution with mean m and variance $v\sigma^2$, and σ^2 has an inverse-gamma distribution with shape parameter α and scale parameter β . With this in mind, for each group G_i we define the hyperparameter vector ϕ_{G_i} as $\langle m_{G_i}, v_{G_i}, \alpha_{G_i}, \beta_{G_i} \rangle$, where each element assumes the corresponding meaning from above. In the next section, we show how this, along with the other parts of the model instantiation, can be used to guide decision making in a multi-agent system.

5.4 Applying the Model

The type of problems we foresee TRAVOS-C being useful for are those which involve an agent choosing rationally which of its peers to interact with. In decision theory, this is generally done by calculating the expected utility of interacting with a particular trustee (see Section 3.1), which can then be compared to other options, such as choosing a different interaction partner or none at all. According to standard theory, and assuming for the moment that we know the value of $\theta_{a_{tr}, a_{te}}$, the expected utility of a_{tr} interacting with a_{te} would be:

$$E[U(O_{a_{tr}, a_{te}}) | \theta_{a_{tr}, a_{te}}] = \int_{-\infty}^{\infty} U(O_{a_{tr}, a_{te}}) p(O_{a_{tr}, a_{te}} | \theta_{a_{tr}, a_{te}}) dO_{a_{tr}, a_{te}} \quad (5.31)$$

where $U(\cdot)$ is some function that specifies the utility of an interaction outcome $O_{a_{tr}, a_{te}}$, from the perspective of a_{tr} . Without knowledge of $\theta_{a_{tr}, a_{te}}$, the best thing to do is marginalise over its value (Equation 5.33), so that the uncertainty surrounding it is fully accounted for. However, this can be a non-trivial problem if $p(\theta_{a_{tr}, a_{te}})$ is hard to evaluate. Alternatively, we can simply estimate $\theta_{a_{tr}, a_{te}}$, either by the mode or the mean of its distribution given the evidence, and then use this to apply Equation 5.31. If the estimation is straightforward, then this is a tractable proposition because, even if there is no analytical solution to $E[U(O_{a_{tr}, a_{te}}) | \theta_{a_{tr}, a_{te}}]$, it only requires integration over a scalar, which can readily be done by numerical integration (Appendix A). However, this only provides an approximation of expected utility, which may be poor if the variance of the

$\theta_{a_{tr},a_{te}}$ distribution is high.

$$E[U(O_{a_{tr},a_{te}})] = E[E[U(O_{a_{tr},a_{te}})|\theta_{a_{tr},a_{te}}]] \quad (5.32)$$

$$= \int_{\Theta^c} E[U(O_{a_{tr},a_{te}})|\theta_{a_{tr},a_{te}}] \cdot p(\theta_{a_{tr},a_{te}}) d\theta_{a_{tr},a_{te}} \quad (5.33)$$

As it stands, in all but the simplest cases, there does not appear to be a closed analytical solution to Equation 5.33, but then neither are there any known closed analytical solutions to estimate $\theta_{a_{tr},a_{te}}$ by its mean or mode. Furthermore, although solutions using infinite series could be sought, it is not obvious that one could be found that would provide an efficiency advantage over numerical integration techniques. For these reasons, we suggest approximating $E[U(O_{a_{tr},a_{te}})]$ directly by calculating Equation 5.33 using numerical integration. For this, we suggest using Monte Carlo techniques (Appendix A) for two reasons:

1. They allow us to estimate any function of $O_{a_{tr},a_{te}}$ (including $E[U(O_{a_{tr},a_{te}})]$) to arbitrary precision, even if the number of unknown parameters is reasonably large.
2. Using them to estimate $E[U(O_{a_{tr},a_{te}})]$ directly does not require significantly more effort than using them to estimate $\theta_{a_{tr},a_{te}}$.

To apply Monte Carlo techniques to this problem, we use them to draw a set of n samples, $\{x_1, \dots, x_n\}$, such that for any function $f(\theta_{a_{tr},a_{te}})$ (e.g. Equation 5.31):

$$\lim_{n \rightarrow \infty} E[f(\theta_{a_{tr},a_{te}})] = \frac{1}{n} \sum_{i=1}^n f(x_i) \quad (5.34)$$

Normally, we could achieve this by sampling from the marginal distribution of $\theta_{a_{tr},a_{te}}$, but this would require the ability to evaluate the density of this distribution, or at least some function proportional to it. Unfortunately, there is no simple way of doing this, but we can achieve the same goal by sampling from the joint distribution for all the parameters in the given instance of the model. This is a more tractable problem because it is relatively easy to find a function proportional to the full joint density. The disadvantage is that this distribution is highly dimensional, depending on the number of agents, which has its own practical issues (Appendix A). On the other hand, one Markov chain will produce estimates for each parameter in the model, so unless we are only interested in a small subset of parameters at a time, we gain some efficiency in this way.

To see how this works in detail, consider a trustor a_{tr} who is interested in the behaviour of n trustees, $a_{te(1)}, \dots, a_{te(n)}$ and consults m reputation sources $a_{rep(1)}, \dots, a_{rep(m)}$. Given our assumptions that both noise and behaviour distributions are normal, this gives the model two sets of parameters:

1. the means and variances of all trustee behaviour distributions; and
2. the means and variances of all reputation noise distributions.

Thus, each sample drawn from the joint distribution is a $l = 2n + 2m$ dimensional vector, of the form:

$$\mathbf{x} = \langle \begin{array}{ccc} \mu_{a_{tr}, a_{te}(1)}, & \dots, & \mu_{a_{tr}, a_{te}(n)}, \\ \sigma_{a_{tr}, a_{te}(1)}^2, & \dots, & \sigma_{a_{tr}, a_{te}(n)}^2, \\ \mu_{\epsilon_{a_{rep}}(1)}, & \dots, & \mu_{\epsilon_{a_{rep}}(m)}, \\ \sigma_{\epsilon_{a_{rep}}(1)}^2, & \dots, & \sigma_{\epsilon_{a_{rep}}(m)}^2 \end{array} \rangle \quad (5.35)$$

By drawing p such samples, $\mathbf{x}_1, \dots, \mathbf{x}_p$, we can approximate the expectation of any arbitrary function $f(\mathbf{x})$ as:

$$E[f(\mathbf{x})] \approx \frac{1}{p} \sum_{j=1}^p f(\mathbf{x}_j) \quad (5.36)$$

provided p is sufficiently large. This includes any function that is dependent on only a subset of parameters, for example the expectation of any particular element of the parameter vector. For instance, if for each i , $\mathbf{x}_i = \langle x_{i,1}, \dots, x_{i,l} \rangle$ and $x_{i,j} \in \mathcal{D}_j$, then from probability theory:

$$E[x_{i,j}] = \int_{\mathcal{D}_j} x_{i,j} p(x_{i,j}) dx_{i,j} \quad (5.37)$$

$$E[x_{i,j}] = \int_{\mathcal{D}_1} \dots \int_{\mathcal{D}_l} x_{i,j} p(x_{i,1}, \dots, x_{i,l}) dx_1 \dots dx_l \quad (5.38)$$

$$E[x_{i,j}] \approx \frac{1}{p} \sum_{i=1}^p x_{i,j} \quad (5.39)$$

which illustrates why approximating expectations using samples from the joint distribution is equally valid as sampling from the marginal distributions. That is, we essentially marginalise over the unused parameters.

To account for group behaviour under this regime is also straightforward, and can be achieved by including the group hyperparameters in the joint distribution. For example, if $a_{te(1)}, \dots, a_{te(n)}$ are divided between g groups labelled G_1, \dots, G_g , then samples from the joint distribution now take the form:

$$\mathbf{x} = \langle \begin{array}{ccc} \mu_{a_{tr}, a_{te}(1)}, & \dots, & \mu_{a_{tr}, a_{te}(n)}, \\ \sigma_{a_{tr}, a_{te}(1)}^2, & \dots, & \sigma_{a_{tr}, a_{te}(n)}^2, \\ \mu_{\epsilon_{a_{rep}}(1)}, & \dots, & \mu_{\epsilon_{a_{rep}}(m)}, \\ \sigma_{\epsilon_{a_{rep}}(1)}^2, & \dots, & \sigma_{\epsilon_{a_{rep}}(m)}^2, \\ m_{G_1}, & \dots, & m_{G_g}, \\ v_{G_1}, & \dots, & v_{G_g}, \\ \alpha_{G_1}, & \dots, & \alpha_{G_g}, \\ \beta_{G_1}, & \dots, & \beta_{G_g} \end{array} \rangle \quad (5.40)$$

set	definition	set	definition	set	definition
\mathcal{S}_{beh}	$\{\sigma_{a_{tr},a_{te}}^2 a_{te} \in \mathcal{T}\}$	\mathcal{S}_{noise}	$\{\sigma_{ca_{rep}}^2 a_{rep} \in \mathcal{R}\}$	\mathcal{S}_{G_i}	$\{\sigma_{a_{tr},a_{te}}^2 a_{te} \in G_i\}$
\mathcal{M}_{beh}	$\{\mu_{a_{tr},a_{te}} a_{te} \in \mathcal{T}\}$	\mathcal{M}_{noise}	$\{\mu_{ca_{rep}} a_{rep} \in \mathcal{R}\}$	\mathcal{M}_{G_i}	$\{\mu_{a_{tr},a_{te}} a_{te} \in G_i\}$
\mathcal{P}_{beh}	$\mathcal{S}_{beh} \cup \mathcal{M}_{beh}$	\mathcal{P}_{noise}	$\mathcal{S}_{noise} \cup \mathcal{M}_{noise}$	\mathcal{P}_{G_i}	$\mathcal{S}_{G_i} \cup \mathcal{M}_{G_i}$
\mathcal{P}_{all}	$\mathcal{P}_{beh} \cup \mathcal{P}_{noise}$	\mathcal{Y}	$\{\phi_{G_i} G_i \in \mathcal{G}\}$	\mathcal{Q}	$\mathcal{P}_{all} \cup \mathcal{Y}$

TABLE 5.1: Parameter set definitions.

The main benefit of this method is that it enables marginalisation over possible group behaviour models (as mentioned in Section 5.2), which will account for any extra information that group behaviour conveys about the likely behaviours of individual agents. A possible secondary advantage is that, if a truster has direct preferences about group behaviour, it can calculate expected utilities that take on board these preferences. For example, an agent may wish to punish groups for erratic behaviour by not interacting with its members. However, we suspect that, in most applications, it would be hard to identify such preferences that cannot be reduced to preferences about individual behaviour.

5.5 A Monte-Carlo Method for TRAVOS-C

Now that we have described the purpose of Monte Carlo methods in TRAVOS-C, we shall describe a particular Monte Carlo process adapted for its application. For this, we use Gibbs sampling (Appendix A) which, while not the most efficient method available, does allow us to get a working model reasonably quickly⁵. Consequently, we do not suggest this mechanism is the best realisation of TRAVOS-C, but only that it shows one possible realisation, which we have used for the empirical evaluation in Section 5.6.

Before detailing this method, however, it is useful to introduce new notation for some commonly used sets. In particular, we introduce five new sets of parameters:

1. the set of noise variances \mathcal{S}_{noise} ;
2. the set of noise means \mathcal{M}_{noise} ;
3. the set of behaviour variances \mathcal{S}_{beh} ;
4. the set of behaviour means \mathcal{M}_{beh} ; and
5. the set of hyperparameters \mathcal{Y} .

Formal definitions for each of these are given in Table 5.1, along with a number of other sets defined in terms of these basic five.

⁵A review of all the techniques used in this section, including alternative techniques that may offer performance advantages, can be found in Appendix A.

Algorithm 2 The TRAVOS-C Gibbs sampler.

Require: $\forall \theta \in \mathcal{Q}$, $\theta^{(0)} \leftarrow$ initial state for θ
 $n \leftarrow$ number of required samples
 $\mathcal{Q}^{(0)} = \{\theta^{(0)} | \theta \in \mathcal{Q}\}$
for $s = 1$ to n **do**
 for $\theta \in \mathcal{Q}$ **do**
 $\theta^{(s)} \leftarrow$ sample from $p(\theta | \mathcal{Q}^{(s-1)} - \{\theta^{(s-1)}\})$
 end for
 $\mathcal{Q}^{(s)} = \{\theta^{(s)} | \theta \in \mathcal{Q}\}$
end for

We can now outline the process in detail. As described in Appendix A, Gibbs sampling produces samples from a high-dimensional distribution, such as the joint parameter distribution in TRAVOS-C, by sampling from a suitable set of conditional distributions that together cover all the dimensions of the target distribution. In our case, this is achieved by drawing independent samples from the conditional distributions of each of the component vectors, $\theta_{a_{tr}, a_{te}}$, ϕ_{G_i} and $\epsilon_{a_{rep}}$ for all a_{te} , a_{rep} and i . The complete sampling mechanism thus follows Algorithm 2.

The independent samples used in this algorithm are generated by a mixture of standard techniques and rejection sampling, depending on the distribution in question. How this is achieved depends on the density of the distribution so, in the following subsections, we outline the sampling mechanisms in each case, by first deriving the conditional densities.

Specifically, we split our discussion into three subsections: Section 5.5.1 derives the equations needed to sample from the hyperparameter distributions required for group behaviour analysis, Section 5.5.2 derives the equations required for sampling from the parameter distributions for individual behaviour, and Section 5.5.3 shows how these can be used to implement the required sampling methods.

5.5.1 Hyperparameter Sampling

Suppose that we wish to draw samples from a distribution with density $p(x)$. To do this, we don't necessarily need the ability to evaluate $p(x)$, as long as we can evaluate some other function, $p^*(x)$, such that $p^*(x) = c \cdot p(x)$, for some (possibly unknown) constant c . When this is the case, we say that $p^*(x)$ is proportional to $p(x)$, and write $p(x) \propto p^*(x)$.

We can now consider the problem of sampling from $p(\mathcal{Q} - \{\phi_{G_i}\} | \phi_{G_i})$, and so need a function that is proportional to it, and can be easily evaluated. To this end, we first note from the model definition that ϕ_{G_i} only directly affects the behaviour distributions of trustees in G_i , and that the behaviour distributions are assumed independent given

their group. As a result, we know that the likelihood $p(\mathcal{Q} - \{\phi_{G_i}\}|\phi_{G_i})$ is equal to:

$$p(\mathcal{Q} - \{\phi_{G_i}\}|\phi_{G_i}) = p(\mathcal{Q} - \{\phi_{G_i}, \mathcal{P}_{G_i}\}|\mathcal{P}_{G_i}, \phi_{G_i})p(\mathcal{P}_{G_i}|\phi_{G_i}) \quad (5.41)$$

$$= p(\mathcal{Q} - \{\phi_{G_i}, \mathcal{P}_{G_i}\}|\mathcal{P}_{G_i})p(\mathcal{P}_{G_i}|\phi_{G_i}) \quad (5.42)$$

which (w.r.t. ϕ_{G_i}) is proportional to:

$$p(\mathcal{Q} - \{\phi_{G_i}\}|\phi_{G_i}) \propto p(\mathcal{P}_{G_i}|\phi_{G_i}) \quad (5.43)$$

$$\propto \prod_{a_{te} \in G_i} p(\theta_{a_{tr}, a_{te}}|\phi_{G_i}) \quad (5.44)$$

Thus, if we assume a uniform prior for ϕ_{G_i} , then from Bayes rule we have:

$$p(\phi_{G_i}|\mathcal{P}_{G_i}) = \frac{p(\phi_{G_i}) \prod_{a_{te} \in G_i} p(\theta_{a_{tr}, a_{te}}|\phi_{G_i})}{p(\mathcal{P}_{G_i})} \quad (5.45)$$

$$(w.r.t. \phi_{G_i}) \quad p(\phi_{G_i}|\mathcal{P}_{G_i}) \propto \prod_{a_{te} \in G_i} p(\theta_{a_{tr}, a_{te}}|\phi_{G_i}) \quad (5.46)$$

which means that to derive each conditional hyperparameter parameter distribution, we need only consider the likelihoods for trustee behaviour distributions belonging to the group in question. Specifically, if we suppose that: $\{\mu_1, \dots, \mu_n\} = \mathcal{M}_{G_i}$, $\{\sigma_1^2, \dots, \sigma_n^2\} = \mathcal{S}_{G_i}$ and $\langle \alpha, \beta, m, v \rangle = \phi_{G_i}$ then from Equations 5.28 to 5.30 we have:

$$p(\phi_{G_i}|\mathcal{P}_{G_i}) \propto \prod_{i=1}^n \frac{\beta^\alpha (\sigma_i^2)^{-\alpha-1}}{\Gamma(\alpha) \sqrt{2\pi v \sigma_i^2}} \exp \left[-\frac{\beta}{\sigma_i^2} - \frac{(\mu_i - m)^2}{2v\sigma_i^2} \right] \quad (5.47)$$

$$\propto \exp \left[-\sum_{i=1}^n \frac{(\mu_i - m)^2}{2v\sigma_i^2} + \frac{\beta}{\sigma_i^2} \right] \frac{\beta^{n\alpha} (\prod_{i=1}^n \sigma_i^2)^{-\alpha-1.5}}{\Gamma(\alpha)^n (2\pi v)^{n/2}} \quad (5.48)$$

Given this, we can consider each component of ϕ_{G_i} individually. To sample from $p(\phi_{G_i}|\mathcal{P}_{G_i})$, we shall prove that:

$$p(\phi_{G_i}) = p(\beta|\alpha)p(\alpha)p(m|v)p(v) \quad (5.49)$$

where $p(m|v)$ is a Gaussian distribution, $p(v)$ is an inverse-gamma distribution, $p(\beta|\alpha)$ is a gamma distribution, and $p(\alpha)$ is a gamcon type II distribution, as defined by Damsleth (1975). There are well known techniques for sampling from each of these types of distribution, apart from the gamcon distribution. So, provided we can specify a sampling method for latter, we can sample from $p(\phi_{G_i})$ by first generating samples from $p(v)$ and $p(\alpha)$, and then using these as known values to sample from $p(\beta|\alpha)$ and $p(m|v)$ respectively.

5.5.1.1 Conditional distribution for β

From probability theory, we know that for any variables x , y and z , $p(x|y, z)$ is proportional (w.r.t. x) to $p(x, y|z)$. Applying this principle to the hyperparameter β , we find from Equation 5.48 that:

$$(w.r.t \beta) \quad p(\beta|\alpha, m, v, \mathcal{P}_{G_i}) \propto p(\beta, \alpha, m, v, |\mathcal{P}_{G_i}) \quad (5.50)$$

$$\propto p(\phi_{G_i}|\mathcal{P}_{G_i}) \quad (5.51)$$

$$\propto \beta^{n\alpha} \cdot \exp\left[-\beta \sum_{i=1}^n \frac{1}{\sigma_i^2}\right] \quad (5.52)$$

Since every p.d.f. must integrate to 1, it follows directly that this is a gamma distribution, and so has the following density with scale parameter ζ and shape parameter k . Furthermore, since this does not depend on m or v , we know that $p(\beta|\alpha, m, v) = p(\beta|\alpha)$.

$$p(\beta|\alpha, \mathcal{P}_{G_i}) = \frac{\zeta^k}{\Gamma(k)} \beta^{k-1} \cdot \exp[-\beta\zeta] \quad (5.53)$$

$$\text{where } k = n\alpha + 1 > 0 \quad (5.54)$$

$$\zeta = \sum_{i=1}^n \frac{1}{\sigma_i^2} > 0 \quad (5.55)$$

5.5.1.2 Conditional distribution for α

At this point we are interested in finding $p(\alpha|m, v, \mathcal{P}_{G_i})$ which, from probability theory, we know is:

$$p(\alpha|m, v, \mathcal{P}_{G_i}) = \frac{p(\alpha, \beta|m, v, \mathcal{P}_{G_i})}{p(\beta|\alpha, m, v, \mathcal{P}_{G_i})} \quad (5.56)$$

We already know $p(\beta|\alpha, m, v, \mathcal{P}_{G_i})$, so all that remains is to find the form of $p(\beta, \alpha|m, v, \mathcal{P}_{G_i})$ and divide by the former. From Equation 5.48 this is:

$$p(\beta, \alpha|m, v, \mathcal{P}_{G_i}) \propto \exp\left[-\beta \sum_{i=1}^n \frac{1}{\sigma_i^2}\right] \frac{\beta^{n\alpha} (\prod_{i=1}^n \sigma_i^2)^{-\alpha-1.5}}{\Gamma(\alpha)^n} \quad (5.57)$$

from which dividing by $p(\beta|\alpha, m, v, \mathcal{P}_{G_i})$ we get:

$$p(\alpha|m, v, \mathcal{P}_{G_i}) \propto \frac{\Gamma(k) (\prod_{i=1}^n \sigma_i^2)^{-\alpha-1.5}}{\zeta^k \Gamma(\alpha)^n} \quad (5.58)$$

$$\propto \left(\prod_{i=1}^n \sigma_i^2\right)^{-\alpha} \left(\sum_{i=1}^n \frac{1}{\sigma_i^2}\right)^{-n\alpha} \Gamma(n\alpha + 1) \Gamma(\alpha)^{-n} \quad (5.59)$$

$$\propto \left[\left(\frac{1}{n} \sum_{i=1}^n \frac{1}{\sigma_i^2}\right) \left(\prod_{i=1}^n \sigma_i^2\right)^{1/n}\right]^{-n\alpha} \Gamma(n\alpha + 1) \Gamma(\alpha)^{-n} n^{-n\alpha} \quad (5.60)$$

$$= C_2 \Gamma(n\alpha + 1) \Gamma(\alpha)^{-n} n^{-n\alpha} \delta^{-n\alpha} \quad (5.61)$$

where C_2 is a normalising constant, and:

$$\delta = \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{\sigma_i^2} \right) \left(\prod_{i=1}^n \sigma_i^2 \right)^{1/n} \quad (5.62)$$

for $\delta > 1$, $\alpha > 0$ and $n > 0$. Here, Equation 5.61 is the p.d.f. for the gamcon type II distribution, as specified by Damsleth (1975). Although the normalising constant, C_2 , is not known to have a closed form, it is clear that the distribution does not depend on m or v , so defines $p(\alpha|\mathcal{P}_{G_i})$.⁶ The constant C_2 does not always exist when $\delta < e$.

5.5.1.3 Conditional distribution of m

In the same vein as β , we find $p(m|\alpha, \beta, v, \mathcal{P}_{G_i})$ by ignoring all constant terms in the joint distribution not depending on m . Thus we have:

$$(w.r.t. m) \quad p(m|\alpha, \beta, v, \mathcal{P}_{G_i}) \propto \exp \left[- \sum_{i=1}^n \frac{(\mu_i - m)^2}{2v\sigma_i^2} \right] \quad (5.63)$$

which from Theorem 5.1 is:

$$p(m|\alpha, \beta, v, \mathcal{P}_{G_i}) = \frac{1}{\sqrt{2\pi\sigma_m^2 v}} \exp \left[- \frac{(\mu_m - m)^2}{2\sigma_m^2 v} \right] \quad (5.64)$$

where

$$\mu_m = \sum_{i=1}^n \frac{\mu_i/v\sigma_i^2}{\sum_{j=1}^n 1/v\sigma_j^2} = \sum_{i=1}^n \frac{\mu_i/\sigma_i^2}{\sum_{j=1}^n 1/\sigma_j^2} \quad (5.65)$$

$$\sigma_m^2 v = \frac{1}{\sum_{i=1}^n 1/v\sigma_i^2} = \frac{v}{\sum_{i=1}^n 1/\sigma_i^2} \quad (5.66)$$

Thus, $p(m|v, \alpha, \beta, \theta_{1:n})$ is a Gaussian distribution, with mean μ_m and variance $\sigma_m^2 v$, and is equivalent to $p(m|v, \mathcal{P}_{G_i})$ since the density does not depend on α or β .

5.5.1.4 Conditional distribution of v

Finally, we find $p(v|\alpha, \beta, \mathcal{P}_{G_i})$ in a similar way to the conditional distribution of α , by first finding the joint distribution of m and v and then dividing by $p(m|\alpha, \beta, v, \mathcal{P}_{G_i})$. So, with respect to v and m we have:

$$p(v, m|\alpha, \beta, \mathcal{P}_{G_i}) \propto v^{-n/2} \exp \left[- \sum_{i=1}^n \frac{(\mu_i - m)^2}{2v\sigma_i^2} \right] \quad (5.67)$$

⁶Since δ is the arithmetic mean of $1/\sigma_{1:n}^2$ divided by its geometric mean, it will always be that case that $\delta \geq 1$. It will only be equal to 1, if all σ_i^2 hold the same value. Nevertheless we may need to legislate for this possibility.

which from Theorem 5.1 is:

$$p(v, m|\alpha, \beta, \mathcal{P}_{G_i}) \propto v^{-n/2} \cdot \exp \left[-\frac{(\mu_m - m)^2 + s^2}{2\sigma_m^2 v} \right] \quad \text{where,} \quad (5.68)$$

where

$$s^2 = \left(\sigma_m^2 \sum_{i=1}^n \frac{\mu_i^2}{\sigma_i^2} \right) - \mu_m^2 \quad (5.69)$$

Now, by dividing Equation 5.68 by $p(m|v, \alpha, \beta, \mathcal{P}_{G_i})$ we find that:

$$p(v|\alpha, \beta, \mathcal{P}_{G_i}) \propto v^{-(n-1)/2} \cdot \exp \left[-\frac{s^2}{2\sigma_m^2 v} \right] \quad (5.70)$$

Then, since Equation 5.70 has the same form as:

$$v^{-\alpha_v - 1} \cdot \exp \left[-\frac{\beta_v}{v} \right] \quad (5.71)$$

where $\alpha_v > 0$ and $\beta_v > 0$, we find that $p(v|\alpha, \beta, \theta_{1:n})$ is an inverse gamma distribution, with parameters defined as follows:

$$\alpha_v = \frac{n-1}{2} - 1 \quad (5.72)$$

$$\beta_v = \frac{s^2}{2\sigma_m^2} \quad (5.73)$$

Moreover, since the density does not depend on α or β , we know that: $p(v|\alpha, \beta, \mathcal{P}_{G_i}) = p(v|\mathcal{P}_{G_i})$.

5.5.2 Parameter Sampling

From the previous section, we have the conditional distributions necessary to draw samples for the group behaviour model. In this section, we turn our attention to the conditional distributions of the parameters for an individual agent's behaviour, which depend on the direct and third party observations available to the truster. Here, our task is simplified by the symmetry that exists between these groups of parameters in the model, in the way they affect the evidence available to the truster. To see why this is the case, let us consider a scenario where a truster is interested in the behaviour distributions of trustees $a_{te(1)}, \dots, a_{te(q)}$, and has evidence from l reputation sources $a_{rep(1)}, \dots, a_{rep(l)}$. If we ignore group-based priors for the moment, and instead assume

uniform priors for all model parameters, the conditional distributions have the form:

$$\begin{aligned} & p(\theta_{a_{tr}, a_{te(j)}} | \mathcal{P}_{all} - \{\theta_{a_{tr}, a_{te(j)}}\}, R_{a_1, a_{te}}, \dots, R_{a_l, a_{te}}) \\ & \propto \prod_{i=1}^m p(R_{a_i, a_{te}} | \theta_{a_{tr}, a_{te(j)}}, \epsilon_{a_{rep(i)}}) \end{aligned} \quad (5.74)$$

$$\begin{aligned} & p(\epsilon_{a_{rep(j)}} | \mathcal{P}_{all} - \{\epsilon_{a_{rep(j)}}\}, R_{a_1, a_{te}}, \dots, R_{a_l, a_{te}}) \\ & \propto \prod_{i=1}^m p(R_{a_i, a_{te}} | \theta_{a_{tr}, a_{te(i)}}, \epsilon_{a_{rep(j)}}) \end{aligned} \quad (5.75)$$

where from Equation 5.12, the likelihood of any particular opinion is given by:

$$\begin{aligned} & p(R_{a_i, a_{te}} | \theta_{a_{tr}, a_{te(i)}}, \epsilon_{a_{rep(j)}}) \propto \\ & \exp \left[-\frac{n_i ((\bar{r}_i - \mu_{a_{tr}, a_{te(j)}}) - \mu_{\epsilon_{a_{rep(i)}}})^2 + s_i^2}{2(\sigma_{a_{tr}, a_{te(j)}}^2 + \sigma_{\epsilon_{a_{rep(i)}}}^2)} \right] (\sigma_{a_{tr}, a_{te(j)}}^2 + \sigma_{\epsilon_{a_{rep(i)}}}^2)^{-\frac{n_i}{2}} \end{aligned} \quad (5.76)$$

where \bar{r}_r is the sample mean, s_i^2 is the sample variance and n_i is the sample size specified by the opinion $R_{a_i, a_{te}}$. From this, it is clear that both the noise and behaviour parameter distributions have precisely the same form, since substituting the noise parameters for behaviour parameters and *vice versa* does not change the equations. This is also true if we introduce direct observations or conjugate priors. In the case of direct observations, these can either be incorporated into the conjugate hyperparameters in the standard way (see Section ??), or be considered a special case of third-party opinions, in which the noise is known to be zero. Likewise, the information given by conjugate priors can also be expressed in the same form as reputation.

To see how this is done, suppose that $\hat{\mu}$ and $\hat{\sigma}^2$ are the mean and variance of interest, and that μ_1, \dots, μ_l and $\sigma_1^2, \dots, \sigma_l^2$ are the additional means and variances present in the likelihood. So, for example, if we are concerned with the conditional distribution of $\theta_{a_{tr}, a_{te}}$, then $\hat{\mu} = \mu_{a_{tr}, a_{te}}$, $\hat{\sigma}^2 = \sigma_{a_{tr}, a_{te}}^2$, $\mu_i = \mu_{\epsilon_{a_{rep(i)}}}$ and $\sigma_i^2 = \sigma_{\epsilon_{a_{rep(i)}}}^2$. Likewise, if we are concerned with a particular noise parameter distribution, then the noise parameters are the parameters of interest, and each pair (μ_i, σ_i^2) denotes the behaviour parameters of a trustee, about which the relevant reputation source has previously expressed an opinion. This gives us the following general form for both noise and behaviour parameter distributions:

$$\begin{aligned} & p(\hat{\mu}, \hat{\sigma}^2 | \mu_1, \dots, \mu_l, \sigma_1^2, \dots, \sigma_l^2) \\ & \propto \exp \left[-\sum_{i=1}^l \frac{n_i ((\bar{r}_i - \hat{\mu} - \mu_i)^2 + s_i^2)}{2(\hat{\sigma}^2 + \sigma_i^2)} \right] \cdot \prod_{i=1}^l (\hat{\sigma}^2 + \sigma_i^2)^{n_i/2} \end{aligned} \quad (5.77)$$

For convenience, it will be useful to talk in terms of the *precision* of interest, which is the reciprocal of $\hat{\sigma}^2$ and is denoted τ . This avoids division by zero when $(\hat{\sigma}^2 + \sigma_i^2) = 0$, which can be useful when implementing the model. If we then introduce an additional

term w_i initially set to $n_i/2$, this gives us the following alternative form:

$$p(\hat{\mu}, \tau | \mu_1, \dots, \mu_l, \sigma_1^2, \dots, \sigma_l^2) \propto \exp \left[- \sum_{i=1}^l \frac{n_i \tau ((\bar{r}_i - \hat{\mu} - \mu_i)^2 + s_i^2)}{2(1 + \tau \sigma_i^2)} \right] \cdot \prod_{i=1}^l \left(\frac{\tau}{1 + \tau \sigma_i^2} \right)^{w_i} \quad (5.78)$$

From this, it is straightforward to see how to incorporate direct observations into this equation, by substituting 0 for μ_i and σ_i^2 , which results in the standard Gaussian likelihood as given by Corollary 5.2. With regard to conjugate priors, such as the group based priors described earlier (Equations 5.28 to 5.30), we know that these have the form:

$$p(\hat{\mu}, \tau | \alpha, \beta, m, v) \propto \tau^{\alpha-1/2} \exp \left[-\tau \beta - \frac{\tau(m - \hat{\mu})^2}{2v} \right] \quad (5.79)$$

So, if we make the following substitutions:

$$\bar{r}_0 = m, \quad \mu_0 = \sigma_0^2 = 0, \quad s_0^2 = 2\beta v, \quad n_0 = \frac{1}{v}, \quad w_0 = \alpha - 1/2 \quad (5.80)$$

we have that Equation 5.79 is equivalent to:

$$p(\hat{\mu}, \tau | \alpha, \beta, m, v) \propto \exp \left[- \frac{n_0 \tau ((\bar{r}_0 - \hat{\mu} - \mu_0)^2 + s_0^2)}{2(1 + \tau \sigma_0^2)} \right] \left(\frac{\tau}{1 + \tau \sigma_0^2} \right)^{w_0} \quad (5.81)$$

which can easily be incorporated in Equation 5.78 without modification. Armed with this information, we can now proceed to derive the conditional behaviour and noise parameter distributions, in terms of $\hat{\mu}$ and $\hat{\sigma}^2$. Specifically, we will derive equations for $p(\hat{\mu} | \hat{\sigma}^2)$ and $p(\hat{\sigma}^2)$. These can then be used to draw independent samples from the joint distribution of $\hat{\mu}$ and $\hat{\sigma}^2$ by first sampling from $p(\hat{\sigma}^2)$ and then using the result to sample from $p(\hat{\mu} | \hat{\sigma}^2)$. In terms of the overall problem, this enables us to sample from the distributions $p(\theta_{a_{tr}, a_{te}} | \mathcal{Q} - \{\theta_{a_{tr}, a_{te}}\})$ and $p(\epsilon_{a_{rep}} | \mathcal{Q} = \{\epsilon_{a_{rep}}\})$ as required by the Gibbs Sampler outlined earlier.

5.5.2.1 Conditional distribution of μ

To derive $p(\hat{\mu} | \hat{\sigma}^2)$, we first note that this is proportional to Equation 5.78, so with respect to $\hat{\mu}$ we have:

$$p(\hat{\mu} | \hat{\sigma}^2, \mu_1, \dots, \mu_l, \sigma_1^2, \dots, \sigma_l^2) \propto \exp \left[- \sum_{i=0}^l \frac{n_i \tau (\hat{\mu} - \bar{r}_i - \mu_i)^2}{2(1 + \tau \sigma_i^2)} \right] \quad (5.82)$$

which from Theorem 5.1 is equivalent to:

$$p(\hat{\mu} | \hat{\sigma}^2, \mu_1, \dots, \mu_l, \sigma_1^2, \dots, \sigma_l^2) \propto \exp \left[- \frac{\tau_r (\hat{\mu} - \mu_r)^2}{2} \right] \quad (5.83)$$

where

$$\tau_r = \sum_{i=0}^l \frac{\tau n_i}{(1 + \tau \sigma_i^2)} \quad (5.84)$$

$$\mu_r = \sum_{i=0}^l \frac{n_i \tau / (1 + \tau \sigma_i^2)}{\sum_{j=0}^l n_j \tau / (1 + \tau \sigma_j^2)} \cdot (\bar{r}_i - \mu_i) \quad (5.85)$$

$$= \sum_{i=0}^l \frac{n_i / (1 + \tau \sigma_i^2)}{\sum_{j=0}^l n_j / (1 + \tau \sigma_j^2)} \cdot (\bar{r}_i - \mu_i) \quad (5.86)$$

$$= \sum_{i=0}^l \frac{n_i \tau (\bar{r}_i - \mu_i)}{\tau_r (1 + \tau \sigma_i^2)} \quad (5.87)$$

$$(5.88)$$

Therefore, the conditional distribution of the mean is a Gaussian distribution, with mean μ_r and precision τ_r . Its p.d.f. is thus given as follows:

$$p(\hat{\mu} | \hat{\sigma}^2, \mu_1, \dots, \mu_l, \sigma_1^2, \dots, \sigma_l^2) = \frac{\sqrt{\tau_r}}{\sqrt{2\pi}} \cdot \exp \left[-\frac{\tau_r (\hat{\mu} - \mu_r)^2}{2} \right] \quad (5.89)$$

5.5.2.2 Conditional distribution of τ

To find the marginal distribution of $\hat{\sigma}^2$, we first divide the joint distribution of $\hat{\mu}$ and $\hat{\sigma}^2$ by the conditional distribution of $\hat{\mu}$, and then remove any superfluous constant terms. From Theorem 5.1 and Equation 5.78, we know that the joint distribution of $\hat{\mu}$ and $\hat{\sigma}^2$ has the form:

$$p(\hat{\mu}, \tau | \mu_1, \dots, \mu_l, \sigma_1^2, \dots, \sigma_l^2) \propto \exp \left[-\frac{\tau_r ((\hat{\mu} - \mu_r)^2 + s_r^2)}{2} - \sum_{i=1}^l \frac{n_i \tau s_i^2}{2(1 + \tau \sigma_i^2)} \right] \cdot \prod_{i=1}^l \left(\frac{\tau}{1 + \tau \sigma_i^2} \right)^{w_i} \quad (5.90)$$

where s_r^2 is defined as:

$$s_r^2 = \sum_{i=0}^l \frac{n_i \tau (\bar{r}_i - \mu_i - \mu_r)^2}{\tau_r (1 + \tau \sigma_i^2)} \quad (5.91)$$

Now, by dividing by Equation 5.89 we get (w.r.t. τ):

$$p(\tau | \mu_1, \dots, \mu_l, \sigma_1^2, \dots, \sigma_l^2) \propto \frac{1}{\sqrt{\tau_r}} \cdot \exp \left[-\frac{\tau_r s_r^2}{2} - \sum_{i=1}^l \frac{n_i \tau s_i^2}{2(1 + \tau \sigma_i^2)} \right] \cdot \prod_{i=1}^l \left(\frac{\tau}{1 + \tau \sigma_i^2} \right)^{w_i} \quad (5.92)$$

$$\propto \frac{1}{\sqrt{\tau_r}} \cdot \exp \left[-\sum_{i=1}^l \frac{n_i \tau (s_i^2 + (\bar{r}_i - \mu_i - \mu_r)^2)}{2(1 + \tau \sigma_i^2)} \right] \cdot \prod_{i=1}^l \left(\frac{\tau}{1 + \tau \sigma_i^2} \right)^{w_i} \quad (5.93)$$

Note that if for all i , $\sigma_i^2 = 0$ then this simplifies as follows:

$$\tau_r = \tau \sum_{i=0}^l n_i \quad \therefore \quad (5.94)$$

$$p(\tau | \mu_1, \dots, \mu_l, \sigma_1^2, \dots, \sigma_l^2) = \frac{\omega^\psi}{\Gamma(\psi)} \tau^{-\psi-1} \exp[-\tau\omega] \quad (5.95)$$

where

$$\omega = \frac{1}{2} \sum_{i=1}^l n_i (s_i^2 + (\bar{r}_i - \mu_i - \mu_r)^2) \quad (5.96)$$

$$\psi = 1/2 - \sum_{i=1}^l w_i \quad (5.97)$$

Thus, when all observations are noise free, τ has a gamma distribution with scale parameter ω and shape parameter ψ .

5.5.3 Sampling Methods for Conditional Parameter Distributions

In the previous sections, we have derived the form of each of the conditional distributions required for Gibbs Sampling. What remains is how to use this knowledge to produce samples from these distributions, which for the most part can be achieved by readily available software libraries.⁷ In particular, sampling solutions for both normal and gamma generated random numbers are given by [Gentle \(1998\)](#), while inverse-gamma samples can be sought by taking the reciprocal of gamma distributed samples.⁸

This leaves only two sets of model parameters that require a specialised mechanism to be devised, namely, the *alpha* hyperparameter distribution (Section 5.5.1.2) and the variance distributions (Section 5.5.1.4). For both of these, we have two choices: either we can generate dependent samples with a technique such as the Metropolis-Hastings method, or independent samples using rejection sampling (Appendix A). Independent samples are always preferable, since estimates based on them converge more quickly. However, there may be a larger overhead in finding a suitable proposal density for rejection sampling, which may outweigh its benefits in certain cases. Therefore, to choose between the two approaches, we need to compare the relative difficulty in finding proposal densities in each case, which depends on the characteristics of the target distributions in our problem.

⁷Libraries that implement suitable sampling methods for gamma and normal distributed variables include the Matlab Statistical Toolbox (<http://www.mathworks.com>) and the NAG Software Libraries (<http://www.nag.co.uk>).

⁸This is the case because the inverse-gamma distribution is defined as the distribution of the reciprocal of a gamma distributed random variable ([Evans et al., 2000](#)).

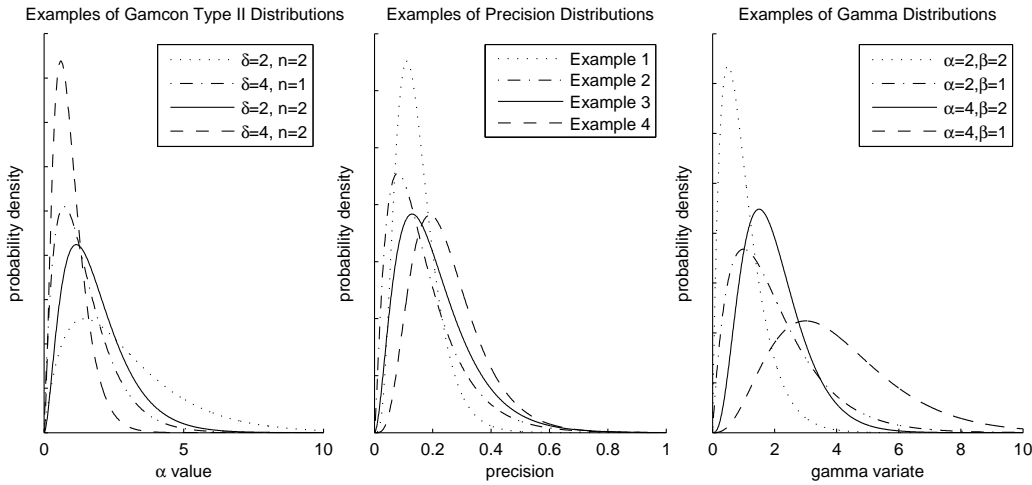


FIGURE 5.4: Examples of gamcon type II, precision and gamma distributions for comparison.

To this end, some examples of both gamcon type II and precision distributions are illustrated in Figure 5.4. As can be seen, both tend to resemble gamma distributions in form, with skewed tails and a single modality. This general shape suggests that, in both cases, a suitable proposal density may be acquired using Laplace’s method (Mackay, 2003). This involves finding the mode,⁹ and then calculating the second log derivative of the distribution, in the region of its mode. The distribution is then approximated with a Gaussian distribution with the same mode, and variance equal to the negative reciprocal of this derivative. That is, for some density $p(x)$, the variance, σ_{approx}^2 , of the Gaussian approximate is given by Equation 5.98, where m is the distribution mode:

$$\sigma_{approx}^2 = - \left(\frac{\partial}{\partial x^2} \ln(p(m)) \right)^{-1} \quad (5.98)$$

$$\sigma_{approx}^2 = \frac{1}{n^2 \phi^{(1)}(nm + 1) - n \phi^{(1)}(m)} \quad (5.99)$$

This technique is used in Garrido (2002) to simulate the gamcon type II distribution using the Metropolis-Hastings method. In this case, the variance of the Gaussian approximate is as shown in Equation 5.99, where $\phi^{(1)}$ is the polygamma function of order 1 (Erdélyi et al., 1953). Garrido points out that the short tails of the Gaussian make it unsuitable to use as a proposal density directly. Instead, a Cauchy distribution is used,¹⁰ by scaling it according to standard deviation of the Gaussian approximate, and shifting it so that it shares the same mode (see Figure 5.5). This gives a similar distribution, but with longer tails that result in better performance.

⁹The mode of a distribution corresponds with the maximum of its p.d.f. or p.f. When we talk about a distribution having multiple modes, we mean that it has multiple local maxima.

¹⁰As with normal and gamma distributions, samples from Cauchy distributions are easily generated using standard libraries and algorithms.

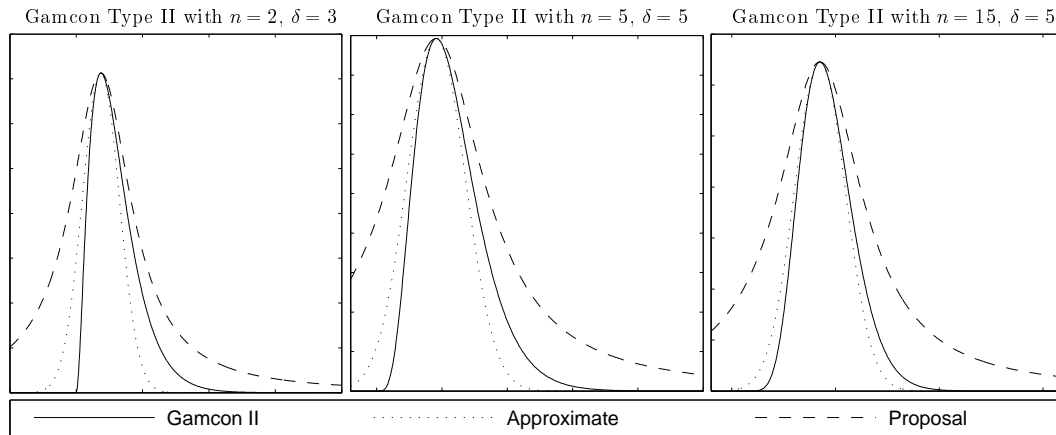


FIGURE 5.5: Example gamcon type II densities with Gaussian and Cauchy approximates.

The main overhead in finding this proposal is the estimation of the mode, which can be achieved using standard hill-climbing algorithms, such as the bisection method, or Quasi-Newton methods (Burden and Faires, 2001). To aid this task, Garrido proves that the gamcon type II distribution is always unimodal (thus, the problem does not suffer from local maxima) and provides the following bounds on its location:

$$\frac{1 - 1/n}{\ln(\delta) + \ln(n/2)} \leq m \leq \frac{2}{\ln(\delta)} \quad (5.100)$$

Given that the shape of the precision distribution is generally similar to the gamcon type II, it seems reasonable that this same technique could also be applied to its simulation. However, if evidence supplied by a truster and its reputation sources differ significantly in the conclusions they support, then it is possible to introduce some local maxima into the precision distribution. As such, algorithms such as simulated annealing (Otten and van Ginneken, 1989), which can deal with local maxima, should be considered as part of a solution.

For both distributions, this technique can also be adapted to provide proposal distributions for rejection sampling. However, in this case we have two extra constraints to consider. First, we need to be able to determine the constant, which is used to multiply the unnormalised target density so that it always lies beneath the proposal density. Second, we need to ensure that the target density fits comfortably within the proposal density to keep the rejection probability reasonably low.

This requirement is illustrated in Figure 5.6. Here, in both parts (a) and (b), the scale of the Cauchy distribution is such that the main region of probable values is narrower than that of the target distribution. This means that the target distribution has to be scaled down low to fit within the proposal density, resulting in a reasonably large

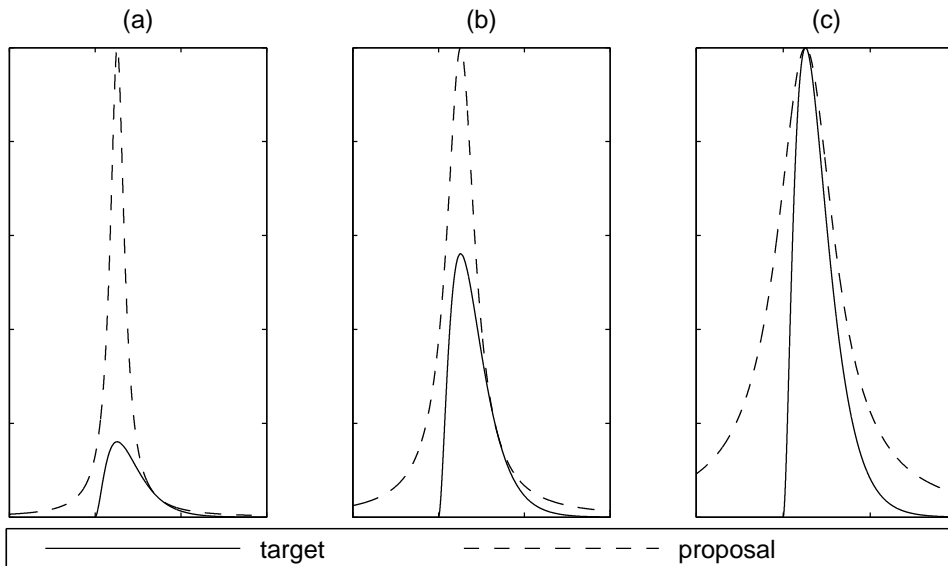


FIGURE 5.6: Cauchy proposal densities for precision distributions.

rejection probability¹¹. Better results are achieved in (c), in which the probable region of the proposal is wide enough to contain the target density, such that the modes of the densities meet. Provided the proposal is not too wide, this can result in better rejection probabilities.

Again, a certain amount of local search may be required to find a good match between the densities. However, we have found that using a scale factor of $2\sigma_{approx}$ gives reasonable results for all the gamcon type II and precision distributions we have tested. Furthermore, in this case the optimal scale factor for the target density tends to be the ratio of the proposal mode over the target mode. That is, if $t(x)$ is the unnormalised target density, $p(x)$ is the proposal density and m is their shared mode, then it is generally the case that:

$$\forall x, \frac{t(x)p(m)}{t(m)} \leq p(x) \quad (5.101)$$

which satisfies the general requirement of rejection sampling. However, there are no guarantees that this is always the case, so some local search is still required to check that the condition is satisfied. Despite this, it seems that, by using these suggested scale factors, a suitable rejection sampling method can be implemented without significant overhead above that of the Metropolis-Hastings method mentioned. Moreover, since independent samples are preferred to dependent samples, this rejection regime should always outperform the Metropolis-Hastings method.

¹¹The rejection probability is proportional to the area of the proposal density outside the region covered by the scaled target density (Appendix A).

5.6 Empirical Study

So far, we have discussed the theoretical aspects of TRAVOS-C, including the definition of the model, and details of how it can be realised and applied using Monte Carlo methods. As described for TRAVOS in Section 4.6, this section evaluates the behaviour of TRAVOS-C in a simulated environment, in which a truster interacts with trustees and reputation sources that behave in certain ways.

More specifically, we present results from a series of experiments in which a truster must assess the reliability of the trustees in its environment, based on varying numbers of direct observations, and reputation with varying degrees of reliability. Together, these results show that TRAVOS-C can not only learn agent behaviour effectively, but that it can do so both when all the assumptions of the model are upheld, and in the presence of certain types of violations of those assumptions.

The rest of this discussion is structured as follows: Section 5.6.1 describes the methodology used to evaluate TRAVOS-C, and collect the results detailed in the preceding sections; Section 5.6.2 presents a number of results concerning the fundamental learning behaviour of TRAVOS-C, given different degrees of information; Section 5.6.3 shows how TRAVOS-C can learn about the reliability of a reputation source, and hence improve its assessment of trustees for which there is little or no direct experience; Section 5.6.4 gives results showing how TRAVOS-C can use information about one reputation source to predict the reliability of other reputation sources; finally, Section 5.6.5 considers cases in which the assumptions made in TRAVOS-C are violated, and shows how TRAVOS-C is robust against certain types of violation.

5.6.1 Experiment Methodology

To evaluate TRAVOS-C, there are three aspects of our experimental methodology that must be considered:

1. the metrics used to assess the model's performance;
2. the simulation process by which the model is tested, and the results recorded; and
3. the methods used to draw conclusions about the model's behaviour and to test for statistical significance.

One way to measure performance is to compare the expected utility calculated by TRAVOS-C to the actual utility a truster receives. This approach is attractive, because it gives an indication of the similarity between decisions made using TRAVOS-C and the optimal decisions a truster could take, given perfect information about its environment.

Algorithm 3 The TRAVOS-C simulation algorithm.

Require: $\mathcal{P} \leftarrow \{a_{te(1)}, \dots, a_{te(p)}\}$ {the set of trustees}

Require: $\mathcal{S} \leftarrow \{a_{rep(1)}, \dots, a_{rep(q)}\}$ {the set of reputation sources}

Require: $\forall i \in \{1, \dots, p\}$, $n_i \leftarrow$ number of direct observations of $a_{te(i)}$

Require: $\forall i \in \{1, \dots, p\}$, $\forall j \in \{1, \dots, q\}$, $m_{i,j} \leftarrow$ number of observations of $a_{te(i)}$ reported by $a_{rep(j)}$
 {Step 1}

for $a_{te(i)} \in \mathcal{P}$ **do**

$dirObs(i) \leftarrow n_i$ samples generated from $a_{te(i)}$ behaviour distribution {direct observations of $a_{te(i)}$ }

for $a_{rep(j)} \in \mathcal{P}$ **do**

$repObs(i, j) \leftarrow m_{i,j}$ samples generated from $a_{rep(j)}$ noise distribution

if $a_{rep(j)}$ is not a liar **then**

$repObs(i, j) \leftarrow repObs(i, j) + m_{i,j}$ samples generated from a_{te} behaviour distribution {observations of $a_{te(i)}$ reported by $a_{rep(j)}$ }

end if

end for

end for

{Step 2}

based on $dirObs$ and $repObs$, estimate all model parameters using Gibbs sampler

{Step 3}

calculate mean absolute errors for each parameter estimate

Unfortunately, an agent's utility depends on its preferences in a particular application domain, which makes it difficult to choose a utility function capable of characterising performance across a range of applications.

For this reason, we choose to measure performance by using TRAVOS-C to estimate the parameters associated with each of the agents in a truster's environment. That is, for each trustee, a_{te} , we estimate the parameters $\mu_{a_{tr}, a_{te}}$ and $\sigma_{a_{tr}, a_{te}}^2$, and for each reputation source, a_{rep} , we estimate the parameters $\mu_{\epsilon a_{rep}}$ and $\sigma_{\epsilon a_{rep}}^2$. The accuracy of these estimates is then measured using the same method employed to evaluate TRAVOS in Section 4.6; that is, by using their mean absolute error using Equation 5.102, in which θ is the parameter being estimated, ϑ , is the estimate, and n is the number of independent simulation episodes used to calculate the mean.

$$\frac{1}{n} \sum_{i=1}^n |\theta - \vartheta| \quad (5.102)$$

Such estimates do not fully determine a truster's ability to make good decisions because the utility of interacting with a trustee may depend on more aspects of its behaviour than the modelling parameters used in TRAVOS-C. However, good estimates of these parameters do imply an ability to characterise at least some aspects of an agent's behaviour, and so provide some indication of how TRAVOS-C should perform in general. With this in mind, we measure the performance of TRAVOS-C in a simulated environment,

in which it is used to estimate trustee behaviour and reputation noise distributions. This is similar to the way in which we evaluated TRAVOS in Section 4.6, except that, in this case, only the trustee behaviour distribution parameter was estimated, since TRAVOS does not explicitly model reputation inaccuracies as added noise. Each experiment consists of a number of independent episodes during which conditions controlled in the experiment remain constant. In turn, these episodes consist of three steps, outlined in Algorithm 3.

First, for each reputation source a_{rep} and trustee a_{te} , a specified number of samples are drawn from a_{te} 's behaviour distribution, to which noise is added, generated from a_{rep} 's noise distribution. In line with the model's assumptions, the resulting noisy samples form the basis for a_{rep} 's opinion about a_{te} , except for certain cases described in Section 5.6.5, in which lying reputation sources are simulated by basing opinions solely on noise, independent of trustee behaviour. Similarly, a specified number of direct observations are generated from each trustee's behaviour distribution, which are made available to the truster without added noise.

Second, given all such direct observations and opinions, Gibbs sampling is used to estimate all trustee behaviour and reputation noise parameters, using the methods described in Sections 5.4 and 5.5. As this is an implementation of Bayesian inference, it requires a prior distribution to be specified for the model parameters. In practice, such a prior could be specified using the group behaviour model (Section 5.2) or some other source of information about the application domain. However, for our purposes, the choice of prior is not significant, as we need only assess TRAVOS-C's performance under different conditions relative to how it performs using only the prior, whatever that prior may be. Thus, we use a normal inverse-gamma distribution (Section 5.3.3) as a prior for the mean and variance of each trustee behaviour distribution, and each reputation noise distribution where, in each case, the hyperparameters used are $\alpha = 2$, $\beta = 10$, $m = 0$ and $v = 100$.

Third, once estimates of each model parameter are collected, we use Equation 5.102 to calculate the mean absolute error for each parameter, over a number of independent episodes, executed under the same experimental conditions. These are then recorded, and used, along with the variance in errors among episodes, to establish the statistical significance of each set of results. Specifically, all claims made in the following sections with regard to these experiments have been tested for statistical significance using analysis of variance techniques. In addition, where results are illustrated using graphs, error bars are displayed using 95% confidence intervals, as is standard practice.¹²

¹²The methods we use to calculate these confidence intervals, along with analysis of variance, are as described by Cohen (1995).

5.6.2 Basic Learning Behaviour

For TRAVOS-C to perform a useful role in evaluating agent performance there are two fundamental hypotheses that should be true, as follows:

1. As a truster gains direct experience of a trustee, its estimation accuracy should improve for both the trustee’s behaviour parameters and the noise parameters of any reputation source that has provided an opinion about that trustee in the past.
2. As the number of observations a reputation source reports for a given trustee increases, the truster’s estimation accuracy should improve for both the trustee’s behaviour parameters and the reputation source’s noise parameters.

To test these hypotheses, we ran a series of experiments during which a truster, a_{tr} , was presented with direct observations of a trustee, a_{te} , along with the opinion of a reputation source, a_{rep} , about a_{te} . During these experiments, no evidence pertaining to any other trustee or reputation source was made available to a_{tr} , and between each experiment, we varied three control variables:

1. the number of direct observations of a_{te} made by a_{tr} ;
2. the number of observations of a_{te} that a_{rep} claims to have made; and
3. the sum of a_{rep} ’s noise variance, $\sigma_{\epsilon a_{rep}}^2$, and a_{te} ’s behaviour variance, $\sigma_{a_{tr}, a_{te}}^2$.

While the first two control variables relate directly to the hypotheses, the third controls the level of difficulty associated with estimation. Standard statistical theory tells us that, to achieve a given estimation accuracy for the parameters of a Gaussian distribution, more observations are required as the variance of the distribution increases¹³ (DeGroot and Schervish, 2002). As TRAVOS-C estimates agent behaviour using observations assumed to be drawn from Gaussian distributions, its performance should not be immune to this effect.

To ensure that the results obtained apply to a general set of agent behaviours, all other aspects of agent behaviour were varied randomly between each episode. Specifically, $\sigma_{a_{tr}, a_{te}}^2$ and $\sigma_{\epsilon a_{rep}}^2$ were determined by assigning a random proportion of their sum to $\sigma_{a_{tr}, a_{te}}^2$, with the remaining proportion assigned to $\sigma_{\epsilon a_{rep}}^2$. This was achieved by generating a random number, P , uniformly distributed on the interval $[0.1, 0.9]$, with which the variances were calculated using Equations 5.103 and 5.104. Given these values, the mean parameters, $\mu_{a_{tr}, a_{te}}$ and $\sigma_{a_{tr}, a_{te}}^2$, were generated from their conditional prior

¹³This is also true for many other useful classes of distribution.

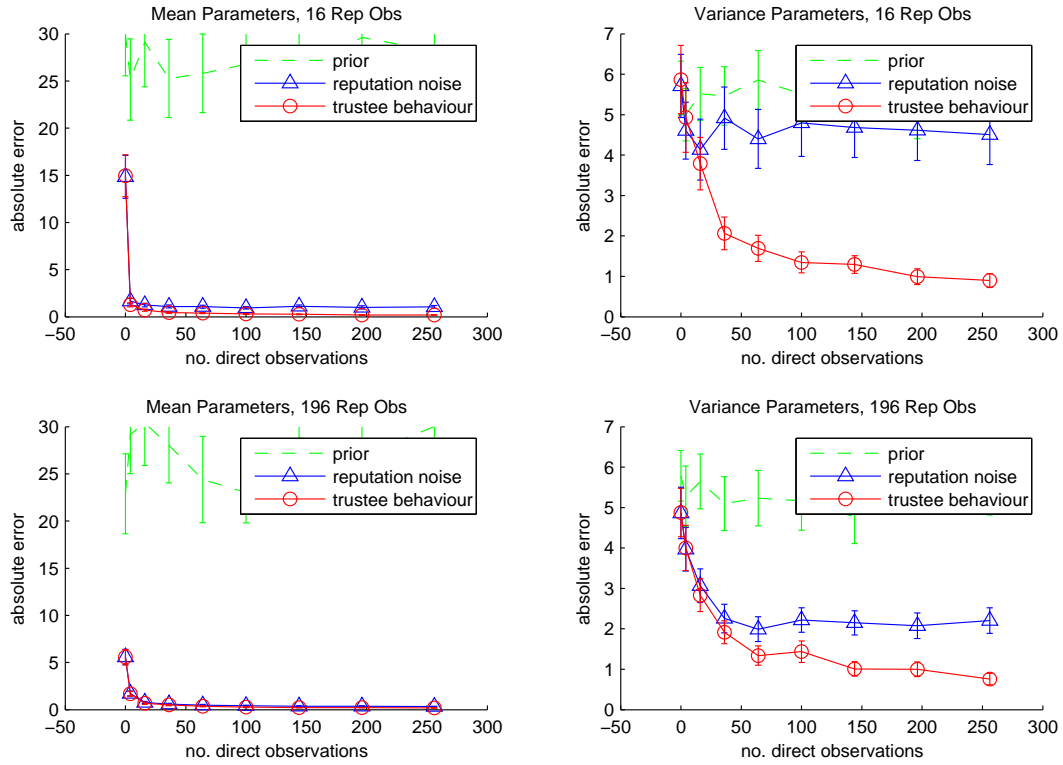


FIGURE 5.7: Parameter estimates with variance sum of 25, varying direct observations.

distributions, as specified in Section 5.6.1.

$$\sigma_{a_{tr}, a_{te}}^2 = P(\sigma_{a_{tr}, a_{te}}^2 + \sigma_{\epsilon_{a_{rep}}}^2) \quad (5.103)$$

$$\sigma_{\epsilon_{a_{rep}}}^2 = (1 - P)(\sigma_{a_{tr}, a_{te}}^2 + \sigma_{\epsilon_{a_{rep}}}^2) \quad (5.104)$$

Selected results from these experiments are illustrated in Figure 5.7, in which the number of direct observations is varied along the horizontal axis of each graph, while the number of reported observations is varied between the top and bottom sets of graphs. In addition, Figure 5.8 gives a similar set of results, except that, in this case, the reported observations vary along the horizontal axes, while the direct observations vary between the top and bottom graphs. In each of these figures, the mean estimation errors achieved for the reputation noise mean, $\mu_{a_{tr}, a_{te}}$, and trustee behaviour mean, $\mu_{a_{tr}, a_{te}}$, are plotted in the left of the figure, while mean estimation errors for the corresponding variance parameters are plotted to the right. For comparison, the estimation error achieved by the model prior is plotted in each graph, showing how the model performs when it has no direct or reported observations.

These results show that, in general, increasing either direct observations or reputation decreases the mean estimation error for each of the model parameters, which is in agreement with the hypotheses stated above. In addition, however, two other notable aspects of behaviour can be observed.

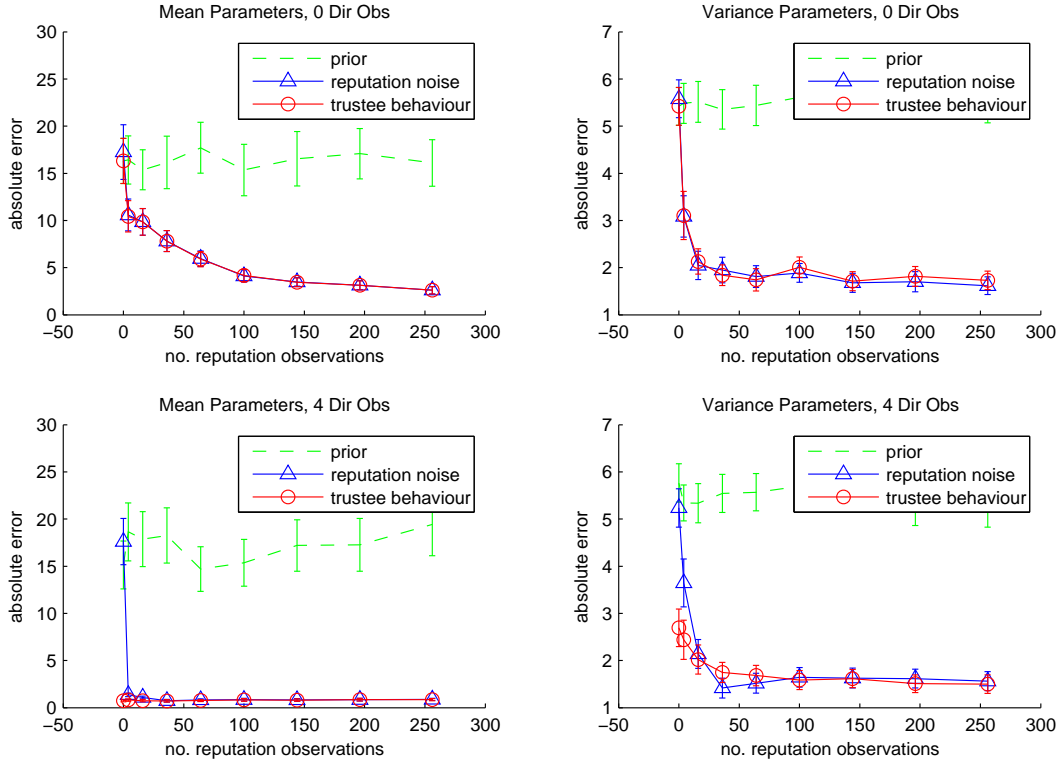


FIGURE 5.8: Parameter estimates with variance sum of 9, varying reputation observations.

First, the number of observations required to significantly improve estimates of the mean parameters is less than required for the variance parameters. This property is intuitive, when we consider the Fisher information (DeGroot and Schervish, 2002) associated with a given number of samples from a Gaussian distribution. Essentially, Fisher information measures the predictive value that each sample has for a given property of its distribution, with higher values indicating more information. In the case of a Gaussian distribution with mean μ and variance σ^2 , the Fisher information of a single sample from that distribution is $1/\sigma^2$ with respect to μ , and $1/2\sigma^4$ with respect to σ^2 . This tells us that it takes significantly more data to obtain the same level of information about the distribution variance, compared to its mean, which is reflected in our results.

Second, as shown in the top two graphs in Figure 5.8, increasing the number of reported reputation observations can improve performance, even when the number of direct observations is 0. This can be attributed to TRAVOS-C learning the sum of the trustee behaviour and reputation noise parameters. That is, although reputation cannot, on its own, be used to determine the noise associated with a reputation source, we can use it to learn the value of the sums $\mu_{atr,ate} + \mu_{earep}$ and $\sigma_{atr,ate}^2 + \sigma_{earep}^2$. These, along with any prior information, can provide some indication of the parameter values, particularly in the case of the variances, because we know that, individually, these must be less than their sum and greater than 0. Moreover, when direct observations are available they illuminate not only a trustee's behaviour, but also the proportion of a reputation source's

opinion that is due to noise.

5.6.3 Learning from Reputation

Although the experiments described in the previous section demonstrate key properties of how TRAVOS-C learns agent behaviour, they do not show how a truster can use reputation to significantly improve estimates of the behaviour of a trustee with whom it has little or no direct experience. To achieve this, a truster should learn from opinions about trustees it has direct experience with, to distinguish reputation sources with low noise, from those with high noise. Using this information, it should then rely significantly on opinions with low noise, while not being misled by inaccurate opinions.

To test for this ability in TRAVOS-C, we placed a truster, a_{tr} , in an environment with two trustees, $a_{te(1)}$ and $a_{te(2)}$, and one reputation source, a_{rep} . During each experiment the number of direct observations was kept constant at 0 for $a_{te(1)}$ and 200 for $a_{te(2)}$, while the number observations reported by a_{rep} for each agent was varied between experiments.

The hypothesis here is that, as the number of observations a_{rep} reports about $a_{te(2)}$ increases, a_{tr} 's estimation accuracy for $a_{te(1)}$'s behaviour parameters also increases when reputation is reliable, or at least stays the same when reputation is relatively noisy. Thus, if this hypothesis is correct, a_{tr} learns about a_{rep} 's reliability from its reported observations of $a_{te(2)}$ (referred to as training observations) and then uses this to determine how to apply a_{rep} 's observations of $a_{te(1)}$ (referred to as test observations).

Similar to conditions described in the previous section, we kept the behaviour variance constant during each experiment (both trustees sharing the same variance), while sampling the behaviour means, independent of each other, from the conditional prior during each episode. With regard to a_{rep} 's noise distribution, two sets of experiments were performed: in one set, $\mu_{\epsilon a_{rep}} = 0$ and $\sigma_{\epsilon a_{rep}}^2 = 0.00001$ (representing near perfect reputation) were used, while in the other set, $\sigma_{\epsilon a_{rep}}^2 = 1000000$ and $\mu_{\epsilon a_{rep}}$ was normally distributed with mean 0 and variance 10000 (representing highly unreliable reputation).

Selected results for experiments with near perfect reputation are illustrated in Figure 5.9, in which the mean estimate errors for $a_{te(1)}$'s behaviour mean are plotted in the top set of graphs, and for the corresponding behaviour variances in the bottom set of graphs. From left to right, the graphs show results for increasing numbers of training observations, with test observations increasing along the horizontal axis. For comparison, the prior estimates are also plotted, as are estimates based on direct observations, equivalent in number to the test observations, with no reputation. The trustee behaviour variance used in this case was $\sigma_{a_{tr}, a_{te(1)}}^2 = \sigma_{a_{tr}, a_{te(2)}}^2 = 81$.

The figure shows that, in general, as a_{tr} learns that a_{rep} 's reputation is essentially noise free, the ability to assess $a_{te(1)}$ based on reputation approaches that based on an

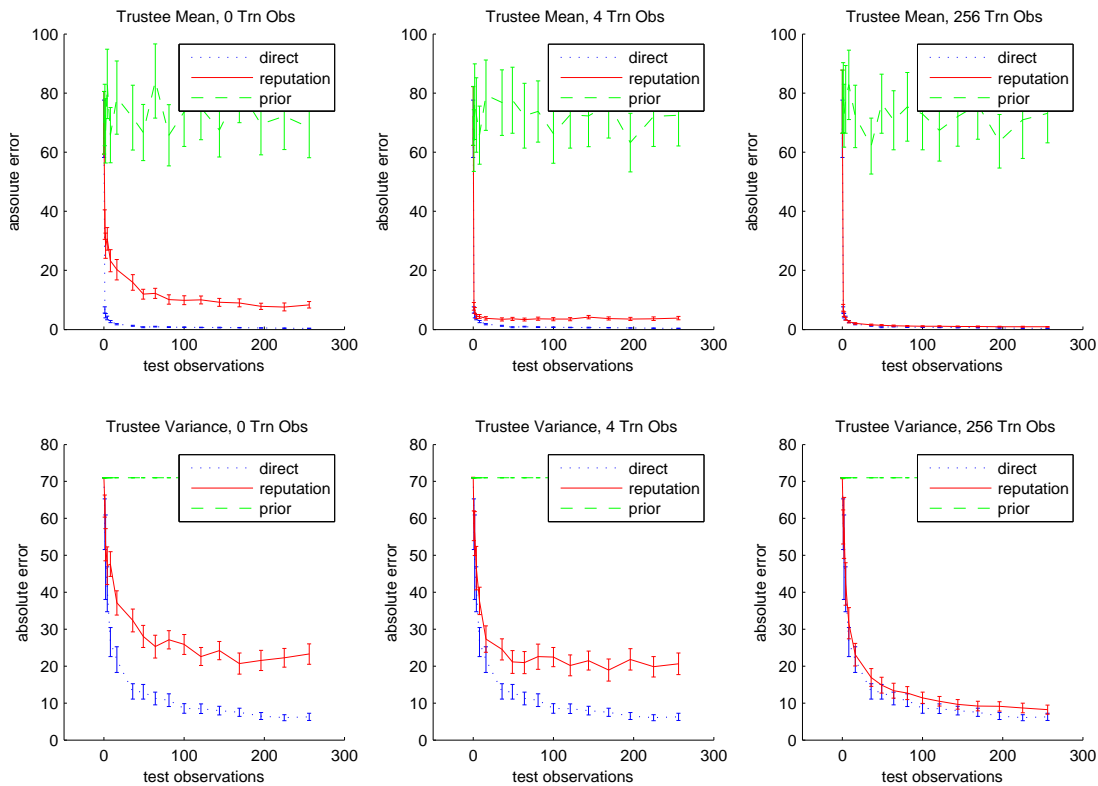


FIGURE 5.9: Behaviour parameter estimates, based on reliable reputation.

equivalent number of direct observations. Moreover, analogous to the results described in the previous section, this convergence takes place more rapidly for the mean than it does for the variance, and even when there are no training observations, some benefit can still be extracted from reputation by using a combination of prior information and parameter sum learning.

In contrast, results for cases in which reputation has significant noise are shown in Figure 5.10, in which the graphs show the estimation errors for $a_{te(1)}$'s behaviour parameters (left) and a_{rep} 's noise parameters (right), plotted against their corresponding prior estimates. As shown in the top set of graphs, there is a spike in estimation error for the behaviour parameters when only one training observation is available with less than five test observations.¹⁴ However, with four or more training observations, this effect disappears, and performance matches that of the prior. This is a positive result because it suggests that, with small amounts of evidence, a truster can learn to ignore the opinion of a reputation source with a significant level of noise.

This effect occurs despite the relatively large estimation errors for the reputation noise parameters, as shown in the figure. These errors can be attributed to the low Fisher

¹⁴This effect is also present when there are no training observations.

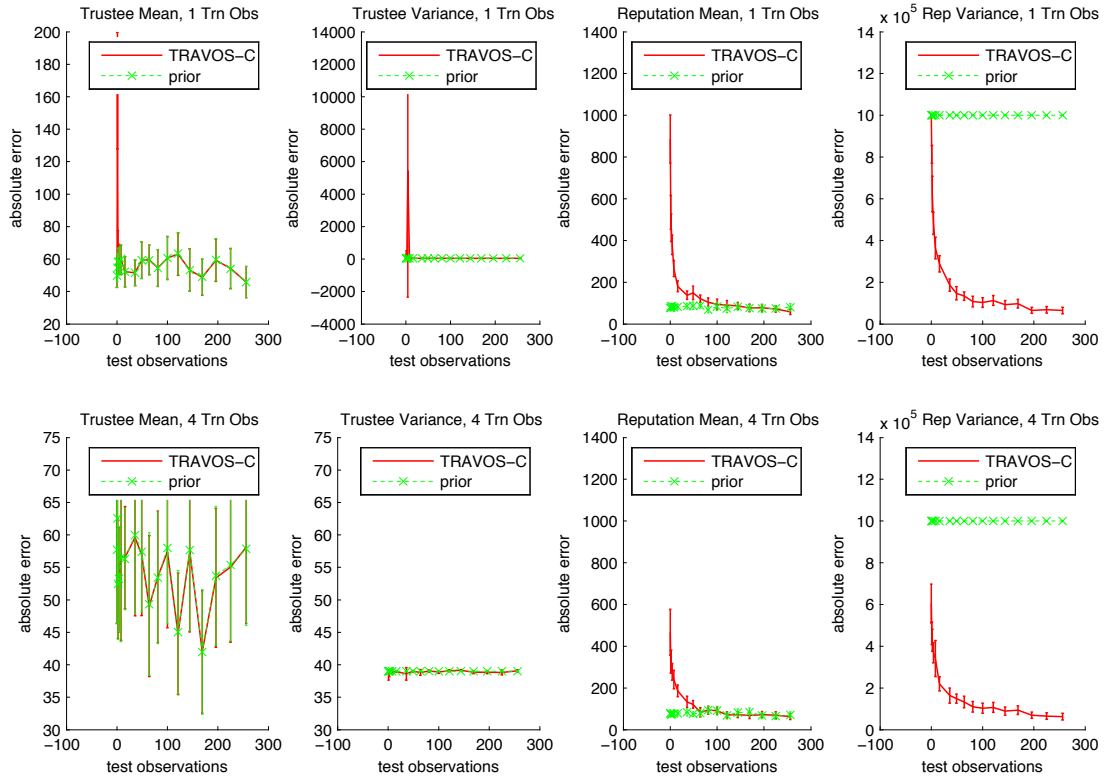


FIGURE 5.10: Behaviour parameter estimates, based on unreliable reputation.

information associated with the reputation, due to the large variance of the noise distribution. However, for a truster to learn to ignore such reputation, it is sufficient to know that the noise variance is high, without necessarily knowing its precise value.

5.6.4 Learning Reputation Source Correlations

So far, the effects on a truster's ability to determine the reliability of a reputation source can mainly be attributed to opinions received from that reputation source, along with direct observations of the trustees those opinions concern. However, at the beginning of the chapter, we claimed that TRAVOS-C can improve its assessment of a reputation source if evidence suggests a correlation between its noise distribution and that of any other reputation source.

For this reason, we evaluated TRAVOS-C in an environment consisting of two reputation sources $a_{rep(1)}$ and $a_{rep(2)}$, and two trustees $a_{te(1)}$ and $a_{te(2)}$. As before, $a_{te(2)}$ was made well known to a_{tr} , with 200 direct observations, while $a_{te(1)}$ had 0 direct observations. However, this time, both reputation sources reported an equal number of observations for $a_{te(1)}$, while only $a_{rep(1)}$ reported any observations for $a_{te(2)}$. In line with previous experiments, the sum of the behaviour and noise parameters was controlled, while the

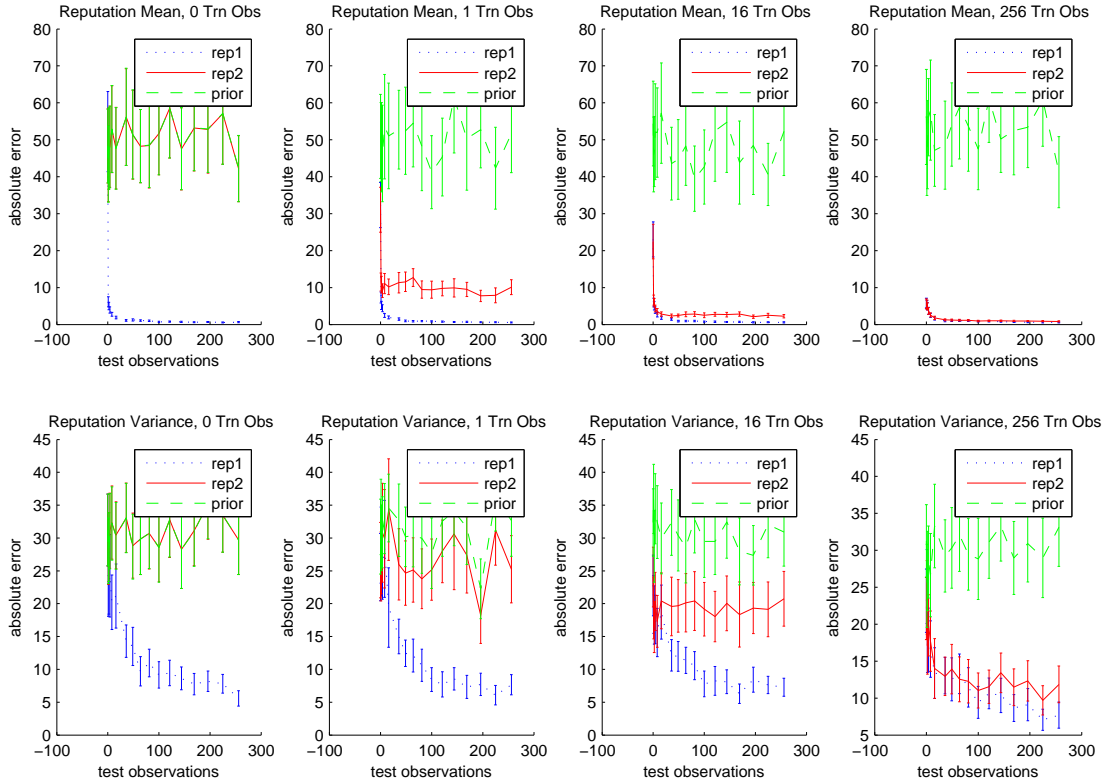


FIGURE 5.11: Reputation noise parameter estimates with evidence for correlation.

proportion of the sum attribute to each variance parameter was allowed to vary, and the mean parameters were drawn from their conditional prior distributions.

Under these conditions, a_{tr} can only directly learn $a_{rep(1)}$'s noise distribution using its direct observations of $a_{te(2)}$. However, even though $a_{rep(2)}$ only reports on $a_{te(1)}$ for which the truster has no direct experience, a_{tr} can still compare the reports of both reputation sources for this trustee. According to the model, any difference between the distributions of these reports must be due to the difference between the noise distributions of each reputation source (since the trustee behaviour distribution is the same). Thus, if a truster were to know this difference, along with the noise distribution of one of the sources, it should be able to infer the distribution of the other.

To validate this claim, we ran experiments in which $a_{rep(1)}$ and $a_{rep(2)}$ had identical noise distributions, varied the number of observations $a_{rep(1)}$ reported for $a_{te(2)}$, and varied the number of observations reported by both sources for $a_{te(1)}$. With respect to these experiments, Figure 5.11 shows the mean estimation errors obtained for the noise distribution parameters of each reputation source, when the total variance shared between noise and behaviour distributions was 162.¹⁵ In each graph, the number of observations reported by $a_{rep(1)}$ for $a_{te(2)}$ varies along the horizontal axis, while the

¹⁵Similar results were obtained using a number of other values for the variance sum.

graphs from left to right plot results for increasing observation numbers reported by both sources for $a_{te(1)}$.

These results show that, as the number of observations $a_{rep(1)}$ and $a_{rep(2)}$ reported for $a_{te(1)}$ increases, a_{tr} 's ability to estimate $a_{rep(2)}$'s parameters, based on knowledge of $a_{te(2)}$, approaches that for $a_{rep(1)}$. In particular, when there are no reported observations for $a_{te(2)}$, there is no evidence to compare $a_{rep(2)}$'s opinions to those of $a_{rep(1)}$. Thus, in this case, direct observations of $a_{te(2)}$ only provide information about $a_{rep(1)}$, and the estimation errors for $a_{rep(2)}$ remain equal to that for the prior.

On the other hand, when there is a significant number of observations from both reputation sources for $a_{te(1)}$, evidence linking the noise distributions of the reputation sources is strong, allowing information about $a_{te(2)}$ to be used to inform predictions about $a_{rep(2)}$.

5.6.5 Performance under Assumption Violations

The results described so far all demonstrate how TRAVOS-C performs when the assumptions of the model are upheld in the environment. However, for TRAVOS-C to have more general applicability, it needs to be robust against at least some types of violation of these assumptions. To evaluate TRAVOS-C under such conditions, we ran three sets of experiments, each representing a different type of violation.

1. We investigated performance when both trustee behaviour and reputation noise distributions are no longer Gaussian, but instead are skewed in one direction. To achieve this, we simulated behaviour and noise using gamma distributions that were transformed such that their mean, variance and skew could be specified as desired. As illustrated in Figure 5.12, the extent to which these distributions violated the Gaussian assumption could be controlled by specifying the skew. In particular, as the skew approaches 0, the resulting distributions have approximately normal morphology.
2. We performed a similar set of experiments using bimodal rather than skewed distributions (Figure 5.13). As with the skewed distributions, these were generated in such a way that the overall mean and variance of the distribution could be controlled, along with the distance between the distributions' modalities. Each distribution was constructed by combining two Gaussian probability density functions, with one Gaussian density per model. The resulting combined p.d.f. is given by Equation 5.105, in which d is a parameter controlling the distance between the modes, and s is chosen to ensure a desired variance for the overall distribution. In the special case where $d = 0$, the resulting distribution is Gaussian with variance

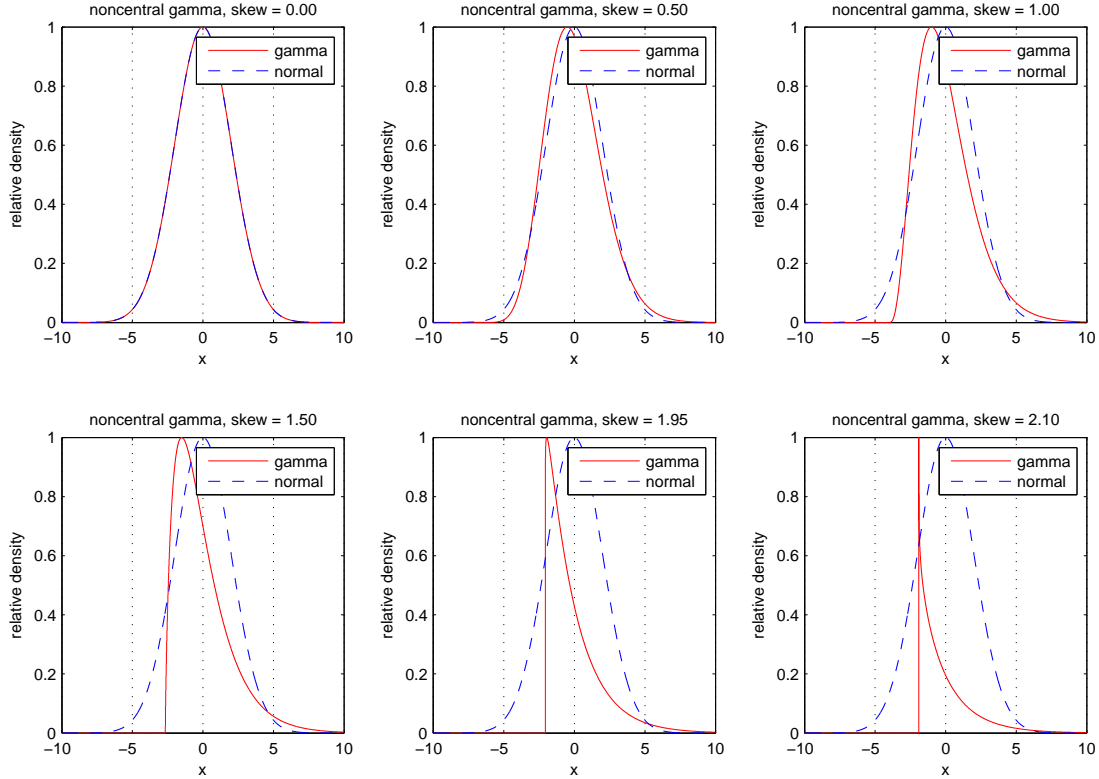


FIGURE 5.12: Example skewed distributions, generated using transformed gamma densities.

s , thus breaking none of the model assumptions.

$$p(x|\mu, s, d) = \frac{1}{\sqrt{8\pi s}} \left(\exp \left[-\frac{(x + \frac{d}{2} - \mu)^2}{2s} \right] + \exp \left[-\frac{(x - \frac{d}{2} - \mu)^2}{2s} \right] \right) \quad (5.105)$$

3. We ran experiments to represent one way in which a reputation source could lie to a truster, in pursuit of its own goals. In this case, rather than adding noise to observations of trustee behaviour, a reputation source's opinion was generated by drawing samples solely from its noise distribution, independent of any true observations of a trustee. This was done to represent cases in which a reputation source simply invents a random opinion, without regard for the trustee's true behaviour.

In the first two sets of experiments, we exposed TRAVOS-C to conditions similar to those described in Section 5.6.2, with 1 trustee and 1 reputation source, and varying numbers of direct and reported observations. However, on these occasions, all behaviour and noise parameters (including the variances) were selected randomly from their prior distributions at the start of each episode, and the way in which these parameters were used in the simulation was changed in line with the particular type of violation being tested. That is, the generated parameters were used to specify the mean and variance

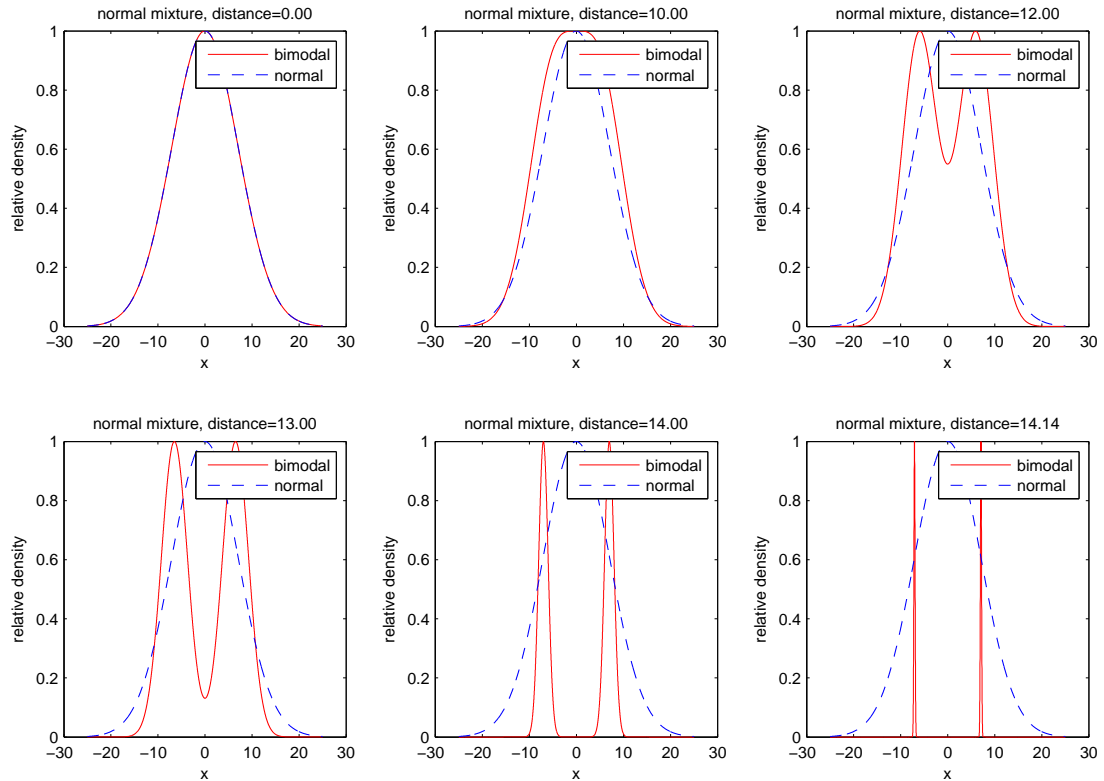


FIGURE 5.13: Example bimodal distributions, generated using a mixture of two Gaussian densities.

of noise and behaviour distributions, which either had a specified skew, or were bimodal with a specified distance between modalities. However, despite using such distributions, we found no significant difference between the results obtained for the modified distributions and those obtained under the model's assumptions.

This suggests that, with regard to estimates of the mean and variance, TRAVOS-C is not sensitive to other properties with respect to the shape of the behaviour and noise distributions. This may be due to the general property that the mean and variance of a sum of random variables is always equal to the sum of the variable means, and the sum of the variable variances, respectively. As stated previously, however, this does not imply that estimates of other distribution properties are just as robust. Thus, the robustness of expected utility calculations may depend on those aspects of a distribution to which a particular utility function is sensitive.

In the final set of experiments, Gaussian distributions were once again used to simulate trustee behaviour and reputation noise, with parameters drawn from their prior distributions between each episode. However, rather than reputation being based on trustee behaviour observations with added noise, opinions from a reputation source were based solely on samples generated from its noise distribution, independent of actual trustee behaviour.

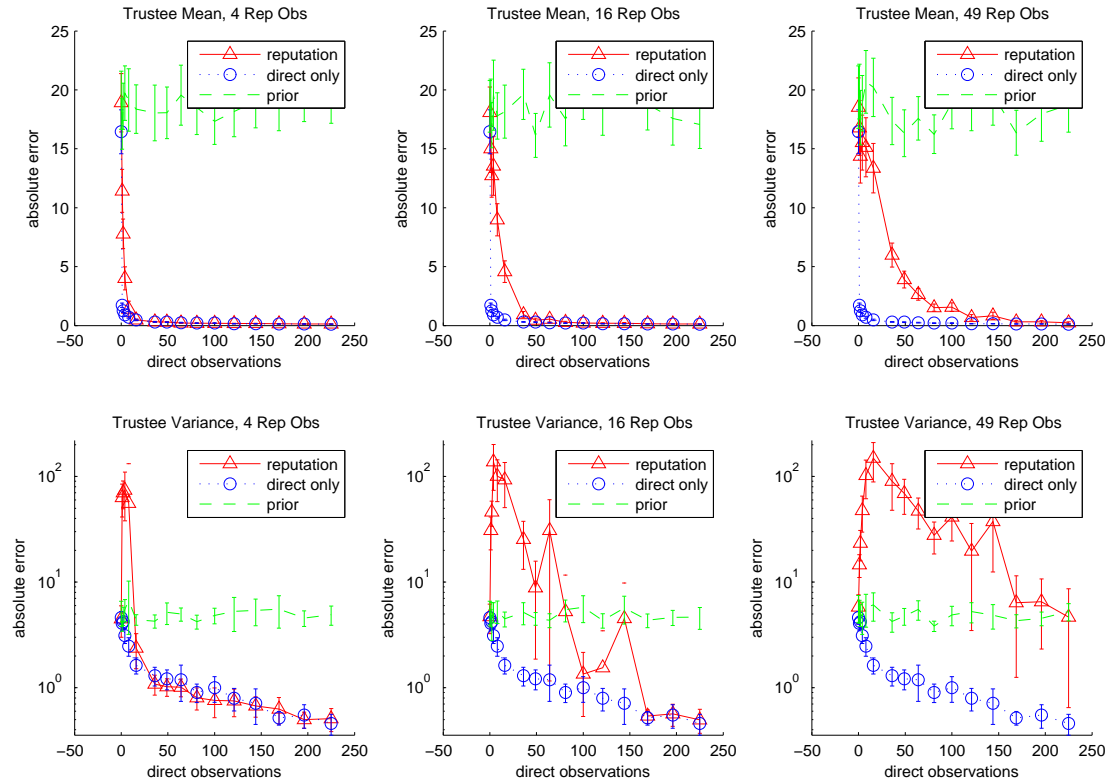


FIGURE 5.14: Behaviour parameter estimates, based on reports from a lying reputation source.

In addition, we varied the number of trustees that the truster was exposed to, each with identical numbers of direct observations, and reported observations from a single reputation source. The reason for this is that, when the reputation source reports for more than one trustee, it provides conflicting evidence about its noise distribution: as the reputation source's reports are generated solely on its noise distribution, reports for different trustees will generally be relatively similar, suggesting that the sum of noise and behaviour parameters should be similar for all agents. Thus, if trustees have dissimilar behaviours, direct observations of each trustee will suggest different values for the noise distribution.

The effect of these experiments on estimates of behaviour distributions is illustrated in Figure 5.14. Here, the number of trustees in the environment is 5, but similar results were obtained from experiments with other numbers of trustees, although the effect was not as bad with just 1 trustee compared to cases with multiple trustees.¹⁶

These results show that, as the number of reported opinions increases relative to the number of direct observations, a truster can be misled by a reputation source with respect to the true behaviour of a trustee. Given enough direct observations, the truster's direct experience can eventually overcome the reputation source's negative impact on

¹⁶In the conditions tested, the number of trustees ranged between 1 and 10.

performance. However, since the reputation source is not tied to report any particular number of observations, it could report having an arbitrarily large observation set to overwhelm any reliable information that a truster may have.

Although this exposes a limitation of TRAVOS-C, one possible solution would be to explicitly build into the model the possibility that reputation is independent of trustee behaviour. However, further investigation would be required to determine how such an approach could deal with a range of possible lying strategies, while maintaining the good performance observed when reputation is reasonably consistent with the current model.

5.7 Summary

In this chapter, we have introduced an extension to TRAVOS, known as TRAVOS-C, which, like its predecessor, fulfills our aim of facilitating an agent in making decisions with regard to its peers. More specifically, TRAVOS-C meets each of the objectives outlined in Section 1.4 by:

- providing a clearly defined mechanism for facilitating rational decision making, through the application of decision theory as described in Chapter 3;
- enabling decisions based on both a truster's direct experience and reputation;
- including mechanisms for efficient communication of reputation, and filtering of inaccurate reputation; and
- enabling reasonable decision making both when a truster is very certain about a trustee's behaviour, and when it has little information about a trustee's behaviour.

In addition, TRAVOS-C has the following three main advantages over its predecessor:

1. TRAVOS-C can assess a trustee based on continuous representations of trustee behaviour, rather than the binary representations used in TRAVOS. Although binary representations may be appropriate in some cases, for example when it only matters that an agent fulfills its obligations and not how it does so, allowing for continuous representations extends the applicability of our work.
2. TRAVOS-C improves on the heuristic reputation filtering mechanism used in TRAVOS by including a new method based on Bayesian analysis. As well as providing a more theoretically sound foundation, this allows a truster to account for more aspects of a reputation source's opinion. In particular, if a group of reputation sources tend to give similar opinions, then information concerning the reliability of one source can be used to assess the reliability of another. Also, even if a reputation source provides significantly biased opinions, these can still count

toward a trustee's assessment, provided the bias is predictable. This is not possible under TRAVOS, because its heuristic method cannot distinguish a predictable bias from one that is unpredictable, providing no information.

3. TRAVOS-C can improve its assessment of a trustee, by considering the behaviour of other similar agents in the system. This is particularly useful when little or no information is available that is specific to the trustee, and this method can adapt its impact on assessment in line with the amount of correlation that exists between agents' behaviour.

To back up these claims, we first defined the model in terms of its assumptions and Bayesian analysis, and showed how it could be applied using decision theory and Monte Carlo methods. We then demonstrated the properties of the system by an empirical analysis. In particular, we showed that TRAVOS-C not only performs well when the model's assumptions are upheld in an agent's environment, but that it is also robust against certain types of violations of those assumptions.

Chapter 6

Conclusions and Future Work

Having described a number of mechanisms for assessing trust and reputation, the purpose of this chapter is to take stock of our work, by outlining both what we have achieved, and what questions remain unanswered. In doing so, the chapter comprises four main sections: Section 6.1 gives an overview of what we have discussed so far, by summarising the main points from each of the previous chapters; Section 6.2 then gives a more detailed view of the contributions we have made to the state-of-the-art; Section 6.3 discusses the main limitations of our work; Section 6.4 discusses the main avenues by which our methods could be extended in future work; and finally, Section 6.5 draws the main conclusions from the thesis.

6.1 Thesis Summary

Trust is a prevalent concept in human society, which is particularly associated with situations in which one entity, a truster, needs to rely on the actions of another entity, known as a trustee. Despite the lack of a single accepted definition, trust can be viewed as the subjective probability with which a trustee will act in a certain way, from the point of view of a truster. This notion of trust is not only important in society at large, but it is also becoming increasingly important in the field of computer science. In particular, we are interested in the role that trust plays in service-oriented systems, such as the Semantic Web and the Grid.

A key objective of these systems is to allow computer resources from different geographical locations, or belonging to different organisations, to be used together seamlessly in support of a common goal. The nature of these systems means that resources from organisations that have competing incentives may be used together, and some resource failure should be expected at anytime. As a result, some researchers have suggested that autonomous software agents, which make decisions without human intervention, could

play an important role in managing resources in such environments. However, if this is to be achieved, these agents must be capable of assessing the trustworthiness of their peers.

In response, our aim was to develop trust assessment mechanisms that could be employed by agents in a service-oriented environment. In particular, we identified two major sources of information that these mechanisms should make use of: (1) the direct experiences of a truster with its peers; and (2) third party experience with a trustee, otherwise known as reputation. However, the amount of information each of these sources provide may vary depending on the situation; in particular, reputation may not always be reliable, due to the view point and incentives of a truster's reputation sources. Thus, even though a truster should make use of these sources, it should be able to deal with inaccurate reputation, and give reasonable results regardless of the amount of reputation available.

Before addressing these aims directly, in Chapter 2 we reviewed existing methods in the literature for solving these and similar problems. Here, we saw that previous models differ both in how they represent trust, and in how they reason about it. In terms of representation, some of the prevailing approaches include the application of Dempster-Shafer theory, probability theory, or more improvised methods. Although each of these may have their place, we believe probability is particularly suited in our context for two main reasons. First, assessing the properties of a system based on past behaviour is one of the fundamental questions that probability theory attempts to answer, and which it achieves through a set of well-established techniques with strong theoretical rationales. Second, probability has a natural interpretation in decision theory, which itself is well suited to facilitating decision making by autonomous agents.

There are three main ways by which existing trust models deal with the inherent lack of reliability in reputation. First, a truster may assume that, out of a group of opinions provided about a trustee, only a minority are likely to be inaccurate. The problem with this approach is that, in many situations, this assumption may be inappropriate. For instance, if no agent has any experience of a trustee, any agent that reports having such experience must be lying, and so is not likely to provide useful information. Second, we may try to discourage lying behaviour by designing systems in which it is always in the best interest of an agent to tell the truth. However, this approach may not always be possible, and in any event, cannot deal with inaccuracies due to reasons other than lying. Finally, we may assess the reliability of a particular reputation source by comparing its opinions to subsequent trustee behaviour: the more correlation we observe between such opinions and behaviour, the more reliable the reputation source can be judged to be. However, among existing probabilistic trust models, this approach has not been addressed in a satisfactory manner.

To begin to address these limitations, Chapter 3 set out a framework by which decision theory and Bayesian analysis can be applied to problems involving trust. In particular, it set out a general approach for making decisions based on a truster's own experiences, and how best to communicate reputation between agents, such that all relevant information is maintained with minimum transmission overhead. Furthermore, to help guide solutions for inaccurate reputation, the chapter categorised the main causes of inaccuracies, along with their effects.

Building on this, Chapter 4 introduced TRAVOS, which instantiates the framework for binary representations of trustee behaviour. TRAVOS includes a mechanism for dealing with inaccurate reputation based on a reputation source's past performance, and has been applied as part of a larger system for managing resources in a service-oriented environment. Finally, Chapter 5 presented TRAVOS-C, which extends the capabilities of TRAVOS in three ways: (1) by using a continuous representation of trustee behaviour, (2) by including an improved Bayesian mechanism for dealing with inaccurate reputation, and (3) by allowing trustees to be assessed based on the behaviour of similar agents in the system.

6.2 Research Contributions

The main contributions of this thesis stem from the specification of the general framework of modelling trust and reputation, and the development of TRAVOS and TRAVOS-C. Together, these show how, by applying standard techniques from statistics and decision theory, an agent can assess the expected benefits of interacting with another agent in a given situation, and so decide which of its peers to interact with, in pursuit of its goals. For example, if an agent has to choose between two providers of a multimedia service, it can use trust to assess how likely each agent is to fulfill its promises, and use trust along with other factors, when making its decisions.

In addition, by applying probability to trust assessment, our methods inherit three key benefits that they share with other probabilistic models of trust. First, by being based on the axioms of probability, these models provide a way of representing beliefs about uncertainty that is consistent and well founded.

Second, by applying well known results, we can derive optimality properties for these mechanisms, under the model assumptions. For instance, decision theory tells us that if an agent has to choose between possible actions, the best choice is always to maximise its expected utility — something which can be directly derived using probability theory.

Third, by applying decision theory along with Bayesian analysis, we can meet our objective of making reasonable decisions regardless of the amount of information available. This is because, through Bayesian analysis, we can calculate the marginal distribution

for a trustee's actions, which accounts for the amount of evidence available in the most appropriate way, given the model assumptions. Then, by applying this in decision theory, a truster can make choices that account for both the risks and potential gains of each choice, given the available evidence.

More significantly, we contribute to the state-of-the-art in three main areas: reputation communication, reputation inaccuracy filtering, and trust assessment based on group behaviour. We elaborate on each of these in the subsections that follow.

6.2.1 Communicating Reputation

As part of our general framework (Chapter 3), we specify a set of guidelines for communicating opinions between agents based on direct experience. These act as a benchmark for transmitting reputation between agents, such that if these are met, all relevant information about an agent's observations are conserved, and this is done with minimum communication overhead. In addition, if these guidelines are adhered to, and an agent's reputation can be assumed to be accurate, then trust models, based on reputation, can be built to reach conclusions that are consistent and as reliable as conclusions based on direct experience.

These guidelines are fulfilled by a number of models, including TRAVOS, which represent trustee behaviour as a binary event, and are based on the Beta Reputation System (BRS). However, TRAVOS-C is the first trust model to address these guidelines for *continuous* representations of trustee behaviour.

6.2.2 Addressing Inaccurate Reputation

In cases where reputation cannot be considered as reliable as direct experience, both TRAVOS and TRAVOS-C implement methods for minimising the impact of inaccurate reputation, each of which has its own separate advantages. The method used in TRAVOS works by comparing the past reports of a reputation source about a trustee, with subsequent direct experience with that trustee. Based on this, it calculates the probability that a trustee's behaviour, on average, lies within a certain margin of error around the reputation source's best estimate. This is then used as part of a heuristic to mediate each source's opinions, such that sources whose opinion accuracy lies outside a margin of error will tend to be ignored completely.

This approach has been shown empirically to outperform the only previous method of its kind (Whitby *et al.*, 2004), which operates on the same BRS derived representation of trust. This is especially important when a significant number of a truster's reputation sources provide inaccurate information, because the method presented by Whitby *et al* assumes that only a minority of opinions are unreliable. Moreover, the method used in

TRAVOS has formed the basis of later work, presented by [Zhang and Cohen \(2006\)](#), which extends the technique to include some of the advantages featured by other existing trust models.

Also building on this, TRAVOS-C presents an improved method of reputation filtering that is derived completely from the assumptions of the model, using Bayesian analysis. In this case, interaction outcomes are represented as real numbers that may, for example, be based on quality of service attributes pertaining to a trustee's performance. Outcomes of interactions between a particular truster and trustee are then assumed to be drawn from a Gaussian distribution with unknown mean and variance. In particular, a truster's direct observations of a trustee are assumed to be drawn from this distribution, while third party observations are assumed to be drawn from the same distribution, but with added Gaussian noise. Each reputation source is associated with a different noise distribution, which the truster may learn through repeated interactions with both trustees and reputation sources.

This approach has several advantages over both TRAVOS and other filtering methods in the literature. First, by applying Bayesian analysis to the model assumptions, we obtain probability distributions for the model parameters, which are provably correct. As such, the model accounts for all evidence and dependencies between parameters that are correct for the model, and can be used to facilitate choices using decision theory in a manner that is theoretically sound, without the need for heuristics. Of course, this does not imply that the assumptions made are correct for every application, but we have shown empirically that TRAVOS-C is robust against many types of violation in its assumptions, and by making its assumptions explicit, it is clear under what conditions the model operates best.

6.2.3 Assessing Trust based on Group Behaviour

To further improve the assessment of a trustee, TRAVOS-C can judge an agent, based on the behaviour of other similar agents in the system. This method is particularly useful for two reasons. First, if neither a truster or its reputation sources have significant experience with a trustee, then assessment based on group behaviour may still provide a significant improvement over assuming no information at all.

Second, it provides a pragmatic solution to the problem of *whitewashing*, in which agents with a poor reputation attempt to improve their standing, by assuming a new identity. In doing so, an agent effectively wipes out any negative information that its peers have about its behaviour, and so is treated just as any other unknown entity in the system. To deal with this, [Zacharia et al. \(1999\)](#) suggest that newcomers to a system should always be assigned the lowest possible rating. However, this may inhibit good market dynamics, by preventing reliable agents from getting a foothold in the market.

As an alternative approach, Sun *et al.* (2005) suggest that newcomers should be judged according to the general behaviour of other newcomers to the system. In doing so, we can adapt our assessments according to the proportion of reliable and unreliable agents that enter the system at any one time.

These advantages can also be claimed for existing models that account for group behaviour, including REGRET (Chapter 2) and Sun *et al.*'s approach. However, these solutions are not directly applicable to probabilistic representations of trust, and require the specification of weights to decide how much impact group behaviour should have.

Our approach is significant in that it automatically adapts to the amount of correlation that exists between the behaviour of a group of agents. That is, only if there is evidence that group behaviour is a strong indicator of an individual agent's behaviour will it have a significant impact on assessment. Conversely, if there a great deal of diversity in the behaviour, then group behaviour will have little or no impact on a truster's assessment.

6.3 Limitations

Although we have a number of methods for assessing the trust that an agent should place in its peers, there are still some open issues for which we do not provide a solution. In particular, we have identified the following three limitations.

Reputation independence assumption In both TRAVOS and TRAVOS-C, we assume that the experiences reported by each reputation source are independent of each other; that is, their observation sets do not intersect. In some cases, it may be desirable to relax this assumption, but this would require some method for either communicating the intersections, or estimating them. Without this, observations that occur in the intersections would have a greater impact on the resulting distribution, leading to an unwanted bias in the results.

Alternative action spaces Through TRAVOS and TRAVOS-C, we provide methods for assessing a trustee when it is appropriate to represent its actions either as binary events, or as real-valued scalars. What we have not offered is a solution to other representations, for example when a truster's preferences depend on multiple dependent attributes, or non-binary discrete action-spaces. To deal with these cases, would require further instantiation of the basic concepts presented in Chapters 3 to 5, but this we consider to be outside the scope of our work.

Overlapping groups With regard to the group behaviour model in TRAVOS-C, we make the assumption that each group that we use for assessment does not intersect. This means, for example, that if we wanted to assess the behaviour of an unknown agent that is blue and comes from Brazil, we could only predict its behaviour based

on other blue agents from Brazil. That is, we could not employ knowledge about agents that are blue and not from Brazil, or from Brazil but not blue. To account for this extra information would require an extra level of complexity, which is not currently present in the model.

6.4 Future Work

In addition to the limitations mentioned above, there are a number of ways in which our mechanisms and their application could be improved, to better aid decision making in a multi-agent system. In particular, we identify the following three areas in which further significant research is warranted.

6.4.1 Dynamic Behaviour

One assumption that we have made in both TRAVOS and TRAVOS-C is that the behaviour of both trustees and reputation sources does not change over time. For many practical applications, this is an unsafe supposition, which we may deal with in one of two ways: either we could assume that agent behaviour does not change significantly over a specified window of time, or we could attempt to model dynamism explicitly in our assessments.

To apply the first of these requires little or no change to our current methods. The only difference would be that, rather than assessing a trustee based on all available observations, we would only base our assessments on observations that have occurred in a certain window of time. The problem with this approach, however, is that it does not specify the duration of the window, nor take account of the precise time in which each observation was made.

A more sophisticated approach would be to model the dynamism in agent behaviour explicitly, by introducing new parameters into our models of trust. This would require a detailed examination of, not only how an agent's behaviour should be modelled under such conditions, but also how reputation is communicated and assessed.

6.4.2 Correlation Between Tasks

Another assumption that we make in our model is that all interactions occur in a similar context. That is, if we wish to know how trustworthy an agent is at providing movie services, we need only consider our past experiences of movie services, and not any other type of service. This is justified because the ability to perform one type of action does not necessarily imply the ability to perform another. Nevertheless, correlations between

an agent's performance for different tasks may exist, and may provide a useful source of evidence. Thus, it may be useful to investigate extensions to our current techniques to take advantage of this knowledge.

6.4.3 Implications of Reputation in Group Learning

An important implication of trust assessment is that a truster will generally choose to interact with agents that, according to the knowledge of the truster, provide better than average performance. Although this seems reasonable, it raises the possibility that a small number of service providers could quickly gain a monopoly position for certain types of service: new agents entering a system may never get a foothold in the market, because no clients will be willing to take a chance on unknown entities.

In human society, this problem is solved by exploration. Although people may generally stick with suppliers that they know, they may occasionally take a risk with a new supplier to judge its performance. In machine learning, such exploration usually falls under the domain of reinforcement learning (Sutton and Barto, 1998), which traditionally considers the problem of individual learners exploring their environment.

Recently, however, research in reinforcement learning has progressed to consider groups of learning entities. Generally, this type of work considers one of two types of problem: (1) agents are self-interested entities, which attempt to learn about each other's behaviour in a competitive environment (e.g. Tran and Cohen (2004)); (2) agents are co-operative members of a team, which attempt to increase group knowledge efficiently by coordinating their actions (e.g. Dutta et al. (2004)). In the former case, agents do not generally share the knowledge that they learn, while in the latter case, agents do share knowledge, but assume that all such knowledge is expressed truthfully.

In our view, agents that share reputation information effectively bridge the gap between these two types of problem. To some extent, trusters are self-interested agents which attempt to learn about the behaviour of other agents to choose the best interaction partners. To achieve this, however, trusters may share information they have about their peers in the form of reputation. This is therefore a cooperative learning problem, with the complication that reputation cannot be assumed to be accurate. Investigation of the use of our current work in combination with reinforcement techniques is therefore warranted.

6.5 Conclusions

Issues of trust are becoming increasingly important in computer science because of the current trend toward large-scale open systems. In particular, we have considered the

case of service-oriented environments, in which resources from different organisations or locations may be used together for a common purpose. Managing such environments can be a challenging problem, because the complement of available resources and the requirements posed upon them may change rapidly over time. As a result, many researchers have suggested that the techniques developed within multi-agent systems should be used to introduce a certain amount of autonomy into the management of such systems.

For such an approach to be effective, however, autonomous agents must be able to assess the trustworthiness of their peers, and make decisions based on these assessments in a clearly justified way. One way to do this is to draw upon the existing techniques of decision theory and statistics, which are already well established approaches for solving these types of problems, and are based on well-founded axioms and reasoning.

However, there are two main factors that make trust assessment problems particularly challenging. First, the amount and type of information a truster has available about a trustee may vary greatly depending on the situation. For example, a truster may have interacted with a trustee many times before, and so have extensive knowledge about the dynamics of the trustee's behaviour. On the other hand, there may be many instances in a large system where a truster comes across a trustee that neither it nor its peers have interacted with before. In both cases, a truster must be able to make reasonable decisions based on the information it has available.

Second, if a truster needs to rely on the third party experiences of its peers, there are two main difficulties that arise that do not arise with a truster's own experience. First, a third party's view of trustee's behaviour may be different from that of the trustee, either because the trustee behaves differently toward different agents, or because the third party does not assess behaviour in precisely the same way as the truster. Second, a truster's reputation sources may have incentives to misrepresent their experiences with a trustee, so as to achieve their own malicious goals. Nevertheless, reputation can be a useful source of information, provided a truster can ward against such possibilities.

To this end, we have developed two models of trust, known as TRAVOS and TRAVOS-C, which can facilitate agent decision making under such conditions. In particular, TRAVOS includes a heuristic mechanism for dealing with inaccurate reputation, and is applicable when an agent's behaviour can be appropriately represented as a binary event. Building on this, TRAVOS-C offers three main advantages over TRAVOS:

1. It reasons about continuous aspects of trustee behaviour.
2. It includes an improved Bayesian mechanism for handling inaccurate reputation.
3. It can assess an agent based on the behaviour of other similar agents in the system.

Both models contribute to the state-of-the-art in computational models of trust, and constitute a significant step toward practical autonomous management of large and open service-oriented environments.

Appendix A

Techniques for Numerical Integration

Often, we need to deal with definite integrals that do not have closed form analytical solutions. Prime examples of this can be found in Bayesian analysis, often occurring when we need to perform expected value calculations, or find the normalising constant for a probability density function. Thus, when analytical solutions are not feasible, we are forced to rely on numerical integration techniques, which can largely be separated into deterministic and non-deterministic methods. Together, these classes cover a large number of methods, including some well established techniques and others that are still the subject of active research. In this appendix, however, we will focus our attention on those techniques that have been employed, or are otherwise appropriate to, the work discussed in this thesis, and in particular a group of non-deterministic techniques known as Monte-Carlo methods. However, to put these in context, we shall begin by giving an overview of some of the main deterministic methods available.

A.1 Deterministic Methods

Many of the most classical deterministic techniques, for example Newton-Cotes rules (Evans and Swartz, 1999), take their inspiration from the definition of definite integration as the limit of a Riemann sum. For example, suppose that a function $f : \mathbb{R} \rightarrow \mathbb{R}$ is continuous on the interval $[a, b]$, and that x_1, \dots, x_n is a sequence of numbers such that:

$$a = x_1 < x_2 < \dots < x_{n-1} < x_n = b$$

In this case, the definite integral of f from a to b can be defined as:

$$\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \sum_{k=1}^{n-1} f(w_k) \Delta x_k \quad (\text{A.1})$$

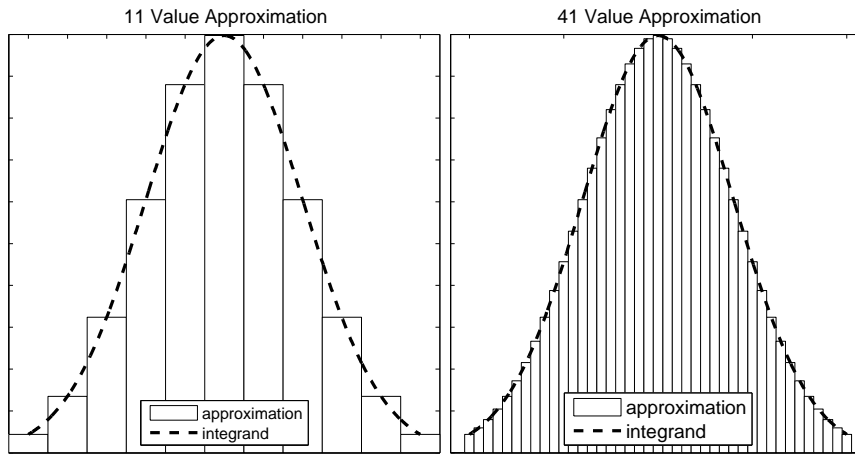


FIGURE A.1: An example of an integrand, approximated with different numbers of rectangles.

where $\Delta x_k = x_{k+1} - x_k$ and w_k is any number in $[x_k, x_{k+1}]$. Essentially, this means that the integral can be approximated by summing the areas of a collection of rectangles that approximate the integrand, as illustrated in Figure A.1. By dividing the domain of integration into an increasing the number of such rectangles, we can arbitrarily increase the accuracy of such approximations at the expense of increasing the size and computational cost of the summation.

More generally, we can partition the domain of integration into one or more regions, and approximate each region by an elementary function which does have an analytical solution to its integral. By increasing the complexity of these approximating functions, we can increase the accuracy of the approximation without increasing the size of summation. However, there is a trade off here, as more complex approximations carry their own computational overhead.

Although such techniques can provide highly accurate results efficiently in low dimensions, in higher dimensions the size of summation required for a given level of accuracy can increase exponentially. Thus, when dimensionality is high, or computational efficiency is paramount, we may need to turn to other methods that do not rely on splitting the integrand into regions.

A simple method for achieving this is Laplace's method ([Evans and Swartz, 1999](#)). In terms of the types of functions it can approximate effectively, this has perhaps narrower applicability than some alternatives, such as variational methods. However, it does provide a simple and effective alternative for certain types of problem. The technique works by approximating a function by an unnormalised Gaussian p.d.f. Based on this, we can integrate the function using the known results for the normal distribution, but the technique also has wider applications; for example, during rejection sampling, which is described below. In one dimension, the method works by basing the variance of the Gaussian approximate on the 2nd log derivative of the integrand in the neighbourhood of

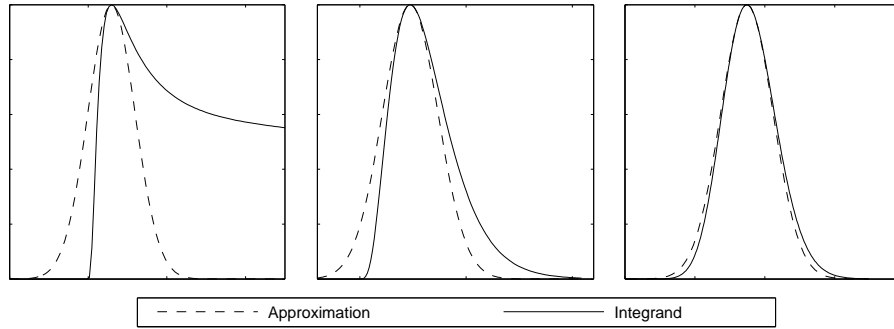


FIGURE A.2: Example Laplacian function approximations.

its maximum. More precisely, if we have an integrand $f(x)$ with an estimated maximum at x_0 , then we approximate $f(x)$ by an unnormalised Gaussian as follows:

$$f(x) \approx f(x_0) \exp \left[-\frac{c(x-x_0)^2}{2} \right] \quad \text{where,} \quad (\text{A.2})$$

$$c = \left. \frac{\partial^2}{\partial x^2} \ln f(x) \right|_{x=x_0} \quad (\text{A.3})$$

For this it can be seen that the approximate is proportional to a Gaussian p.d.f. with variance $\frac{1}{c}$. As such, we know from the normalising constant of this Gaussian that:

$$\int_{-\infty}^{\infty} f(x) dx \approx f(x_0) \sqrt{\frac{2\pi}{c}} \quad (\text{A.4})$$

Beyond this, the technique can be generalised to higher dimensions by calculating the Hessian matrix¹ of the integrand's natural logarithm. This can then be used to approximate the integrand by a multidimensional Gaussian with appropriate covariance matrix. However, as is clear from this description, the Laplace method is only applicable if the integrand can be approximated effectively by a Gaussian density function. Sometimes this can be achieved by applying some transformation to the integrand, but, in any event, the function should have one predominant maximum, with relatively symmetric sloping sides that decreases monotonically toward zero. To illustrate this, Figure A.2 shows some example functions with corresponding Laplacian approximates.

A.2 Monte Carlo Methods

So far, the numerical integration techniques that we have described suffer from one of two problems: either they place a hard limit on the level of accuracy they can achieve for certain types of functions, or they quickly become intractable as dimensionality grows.

¹The Hessian matrix is the multidimensional equivalent of the 2nd derivative in one dimension. It contains entries of all the partial 2nd derivatives of a function (Khuri, 2003).

In contrast, Monte Carlo techniques occupy the middle ground for two main reasons: (1) they are anytime algorithms, in that greater accuracy can always be achieved by spending more compute time on estimation; and (2) they are less sensitive (though not always immune) to the effects of dimensionality, compared to many deterministic techniques.

The defining property from which Monte Carlo techniques derive their strength is that they simulate a stochastic process. They are in fact sampling methods, which can be used to draw samples from almost any probability distribution of choice. Integration problems are solved as a by product of this, by reformalising a definite integral as an expected value calculation. This makes expected utility calculations a particularly natural application of Monte Carlo techniques, but other integration problems can also be solved in this way, even if an expected value is not their aim.

To give an example of how these work, suppose we wish to perform definite integration of a vector x over domain \mathcal{X} , and we can express this integral as the expected value of some function $f(x)$ with p.d.f. $p(x)$ (Equation A.5). Then, according to probability theory, we can approximate this integral by averaging the value of $f(x)$ over n i.i.d samples, $\{x_1, \dots, x_n\}$ drawn from $p(x)$ (Equation A.6).

$$E[f(x)] = \int_{\mathcal{X}} f(x)p(x) dx \quad (\text{A.5})$$

$$E[f(x)] \approx \sum_{i=1}^n f(x_i) \quad (\text{A.6})$$

This results in an anytime algorithm because we can perform the calculation with any reasonable value of n . However, the larger we allow n to be, the larger the expected accuracy will be. Furthermore, the accuracy of this technique is not directly dependent on the dimensionality of the domain, but only on the variance of the distribution specified by $p(x)$. However, as we shall see, dimensionality can cause other problems, because it can increase the difficulty associated with drawing independent (or effectively independent) samples.

With this mind, we shall now explore some of the mainstream techniques for sampling from a distribution. In the main, there are two ways to do this: either we attempt to draw independent samples from the distribution, or we can draw dependent samples using a Markov chain. Independent samples are always preferred, because generally they provide a higher level of accuracy with a lower number of samples. However, generating independent samples is not always easy, in some cases carrying a higher computational overhead than they are worth. On the other hand, Markov chain Monte Carlo (MCMC) techniques work by generating a sequence of samples in which each sample depends on the sample that directly precedes it in the chain. Provided the marginal distribution of each sample (with all other samples unknown) is the distribution we wish to sample from, then the approximation will still converge on the desired result. However, the

larger the correlation between each sample and its predecessors, the more samples it will take to achieve a desired level of accuracy.

A.3 Independent Samples and Rejection Sampling

To generate truly independent samples from a given distribution, we have two main possibilities: inversion methods, or rejection sampling (Mackay, 2003). Both of these techniques assume that we can generate samples from some elementary distribution, such as a uniform distribution, for which there are many existing algorithms (Gentle, 1998). The samples are then manipulated in some way, so that the distribution of the modified samples is as desired.

In the case of inversion methods, this is achieved by generating uniform samples on the interval $[0, 1]$ and then transforming these using the inverse distribution function of the target distribution. However, in many cases, calculating the inverse distribution function is not a tractable proposition. Therefore, unless the peculiarities of the problem allow for some other sampling regime, we must turn to rejection sampling.

The main difference between rejection sampling and sample transformation, is that rejection sampling does not keep all of the samples it actually generates. Instead, samples are first proposed, and then selectively thrown away, such that the remaining samples have the desired distribution. To see how this is achieved, suppose that we have two distributions: (1) a *target* distribution with p.d.f., $P(x)$, which we wish to sample from; and (2) a *proposal* distribution with p.d.f., $Q(x)$, which we actually *can* sample from. We then need to be able to evaluate functions $P^*(x)$ and $Q^*(x)$ that are proportional to $P(x)$ and $Q(x)$ respectively. That is, it must be the case that:

$$[\forall x \in \mathcal{X}, P(x) = aP^*(x)] \wedge [\forall x \in \mathcal{X}, Q(x) = bQ^*(x)]$$

where a and b are two (possibly unknown) constants and \mathcal{X} is the domain of x . In addition, it must also be the case that:

$$\forall x \in \mathcal{X}, Q^*(x) > P^*(x)$$

Although the last of these constraints may not always be easy to fulfill, that we only need of evaluating functions proportional to the two densities gives us significant freedom. In particular, it releases us from the necessity to evaluate the normalising constant for the target distribution, if this is not of direct interest. For example, in Section ??, the expected value calculations given in Equations 3.3 and 3.4, both involve the posterior density $p(\theta|X)$ from Equation 3.6. In this case, we could use $P^*(x) = p(\theta)p(X|\theta)$, so avoiding the need to evaluate $p(X)$, which we have already stated can be problematic.

Algorithm 4 The rejection sampling algorithm.

```

for  $s = 1$  to  $n$  do
   $accepted \leftarrow false$ 
  while  $accepted$  is  $false$  do
     $s \leftarrow$  sample from  $Q$ 
     $u \leftarrow$  uniform sample from interval  $[0, Q^*(s)]$ 
    if  $u \leq P^*(s)$  then
       $x_i \leftarrow s$ 
       $accepted \leftarrow true$ 
    end if
  end while
end for

```

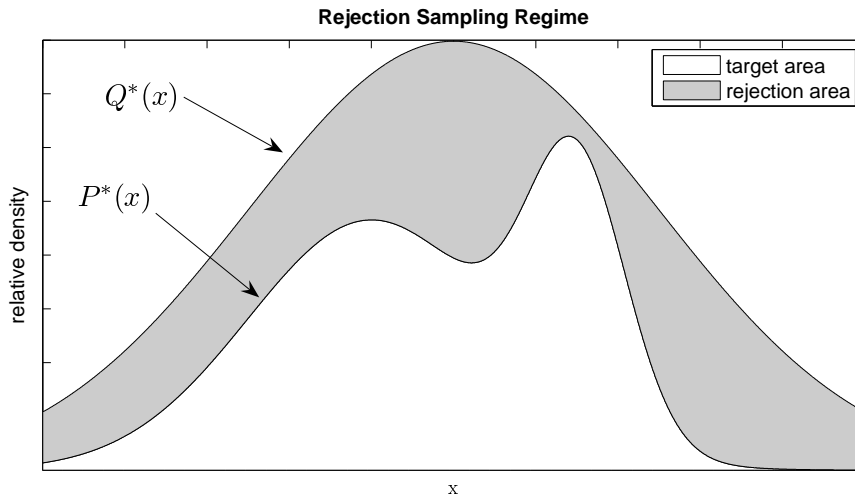


FIGURE A.3: Example rejection sampling regime.

With this in mind, the method proceeds according to Algorithm 4. Here, the first step is to draw a sample s from the proposal distribution, following by a second sample u , drawn uniformly from the interval $[0, Q^*(s)]$. Based on this, s is rejected if $u > P^*(s)$ and kept otherwise. The reason this works is illustrated in Figure A.3. In effect, the vector $\langle s, u \rangle$ is uniformly drawn from the area under the curve of $Q^*(x)$. Since any samples that are drawn in the shaded area between $Q^*(x)$ and $P^*(x)$ are rejected, the remaining vectors are uniform samples from the area under $P^*(x)$. As a direct consequence of this, the s components are distributed according to the $P(x)$ which is the desired result.

Although this is true for any proposal density that satisfies the stated constraints, not all such densities are equal. The distinguishing factor is the rejection probability, which marginalised over x , is given as:

$$\int_{\mathcal{X}} \frac{P^*(x)}{Q^*(x)} Q(x) dx \quad (\text{A.7})$$

Normally, generating samples from the proposal density is relatively inexpensive. However, the number of proposal samples required per accepted sample, grows proportionally

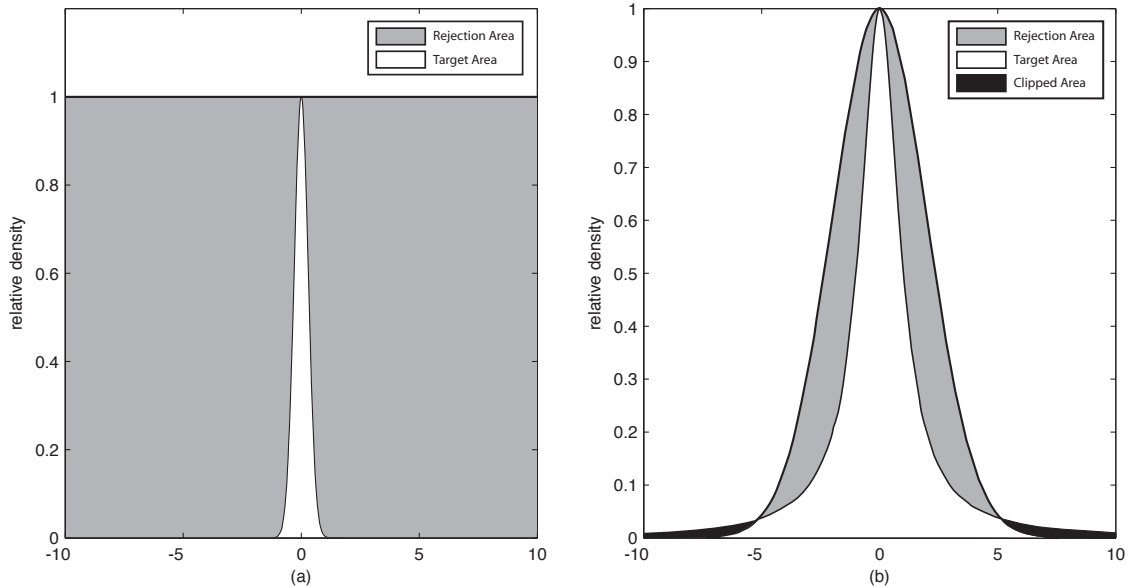


FIGURE A.4: Examples of poor rejection regimes.

with the rejection probability. Thus, although it may be acceptable to reject 50%, or even 95% of samples, rejecting 99.99% of samples may prove to be inhibitingly expensive, depending on the application. A certain amount of optimisation on the part of the proposal density may thus be desirable, although achieving near zero rejections is unlikely to be necessary. However, to achieve even a 1% acceptance rate, the proposal density needs to strike a balance between adequately covering the full support of the target distribution, and expending too much density in areas where the target has little or no density.

These two extremes are exemplified in Figure A.4. Here, part (a) illustrates a very narrow target distribution with a uniform proposal distribution covering a much wider domain. The rejection probability in this case will be high because there will be a relatively low probability of generating samples in the region of the target. On the other hand, part (b) shows a standard normal distribution being used to propose samples for a Cauchy distribution of similar scale. This case is actually a misnomer, because it is impossible to use a Gaussian to adequately propose samples for a Cauchy, due to the Cauchy distribution's wider tails. In effect, the area shaded black in the figure would be clipped from the target distribution, and so any accepted samples would actually be distributed according to the remaining area of the Cauchy, shaded white.

The search for viable proposal densities is also hindered by dimensionality. This is because, as dimensionality grows, it becomes more and more difficult to find $Q^*(x)$ always greater than $P^*(x)$, without resulting in impractical rejection probabilities. Moreover, it is generally the case that the acceptance rate becomes exponentially small as dimensionality increases. To some extent, this can be overcome if the distribution can be

Algorithm 5 The Metropolis-Hastings algorithm.

```

 $x_0 \leftarrow$  initial state
for  $s = 1$  to  $n$  do
   $x' \leftarrow$  sample from  $Q(x|x_{i-1})$ 
   $a = (P^*(x')Q(x_{i-1}|x')) / (P^*(x_{i-1})Q(x'|x_{i-1}))$ 
  if  $a \geq 1$  then
     $x_i = x'$ 
  else
     $x_i = x_{i-1}$ 
  end if
end for

```

broken into its low dimensional conditionals. For example, suppose we wish to sample the joint distribution of three scalars, with p.d.f. $P(x_1, x_2, x_3)$. This can be achieved by sampling first from $P(x_1)$ followed by $P(x_2|x_1)$ and then $P(x_3|x_1, x_2)$. In each case, the samples from the current distribution supply the conditionals for each subsequent distribution. This, however, is not always feasible, so rejection sampling is not considered an appropriate technique for much more than one dimensional problems.

A.4 The Metropolis-Hastings Algorithm

When truly independent sampling can not be efficiently achieved, then turning to dependent sampling may offer a viable alternative. One of the simplest of such approaches is the Metropolis-Hastings algorithm (Mackay, 2003). This works in a similar way to rejection sampling, except that the proposal density need not be similar to the target density for the technique to be effective. Instead, samples are proposed from a density $Q(x|x')$, which is *dependent* on the current state x' . That is, instead of generating a set of independent samples, we now generate a sequence of samples, x_1, \dots, x_n , in which each x_i is directly dependent on its predecessor, x_{i-1} , starting from some initial state x_0 . Such a sequence is known as a Markov chain, which is why Metropolis-Hastings and other sampling regimes that use this technique are known as Markov Chain Monte Carlo (MCMC) methods.

Although in using a Markov chain, we lose the ability to generate independent samples, the advantage over rejection sampling is that $Q(x|x')$ no longer needs to be similar in form to $P(x)$, something which can be difficult to achieve for complex distributions. Instead, the method proceeds according to Algorithm 5, in which each sample x_i is proposed from $Q(x_i|x_{i-1})$ and accepted if:

$$\frac{P^*(x')Q(x_{i-1}|x')}{P^*(x_{i-1})Q(x'|x_{i-1})} \geq 1 \quad (\text{A.8})$$

Algorithm 6 The Gibbs sampling algorithm.

```

 $\mathbf{x}^{(0)} = \langle x_1^{(0)}, \dots, x_m^{(0)} \rangle$  is initial state
for  $i = 1$  to  $n$  do
  for  $j = 1$  to  $m$  do
     $x_j^{(i)} \leftarrow$  sample from  $P(x_j | \{x_l^{(i-1)}\}_{l \neq j})$ 
  end for
end for

```

If this is not the case, the proposed sample is rejected. However, rather than leaving the sequence untouched, rejecting a sample causes the previous sample state to be written again to the sequence.

The key property of this process is that the distribution of x_i tends to $P(x)$ as $i \rightarrow \infty$ and, as a result, it can be shown that samples generated by Metropolis-Hastings can be used to estimate expected values through Equation A.10, where \mathcal{D} is the domain of x . However, how fast this convergence takes place is another matter, which determines how many samples are required to achieve a given level of accuracy, relative to a given number of independent samples.

$$E[f(x)] = \int_{\mathcal{D}} f(x)P(x) dx \quad (\text{A.9})$$

$$\approx \frac{1}{n} \sum_{i=1}^n x_i P(x_i) \quad (\text{A.10})$$

As with the acceptance rate in rejection sampling, this convergence rate is determined by the size and shape of the proposal distribution, relative to the target distribution. If the proposal density is too large compared to the target, then many samples will be proposed that have low probability according to the target, and so will be unlikely to be accepted. On the otherhand, if the proposal density is too small, then the sample sequence will exhibit a random walk behaviour, which will take a long time to adequately explore the target density. As well as this, it can be difficult to assess how many samples are required to achieve a certain level of accuracy, although many techniques have been proposed to tackle this problem to some extent (Evans and Swartz, 1999).

Once again, dimensionality plays a role in how difficult it is to achieve an efficient sampling mechanism. This time, however, the relationship between the number of dimensions and the size of simulation required for a given level of accuracy is essentially quadratic, rather than exponential. Thus solutions involving a reasonably high number of dimensions are at least possible, though we may still need to run long simulations to achieve them.

A.5 Gibbs Sampling

Gibbs Sampling (Mackay, 2003) is another example of a MCMC technique, which can be applied to high dimensions. In fact, Gibbs sampling can be shown to be a specialise case of Metropolis-Hastings, so any statement that is true for Metropolis-Hastings methods is also true of Gibbs Sampling. What disguises it from the general case is that it does not involve any adjustable parameters in the form of a proposal density, so it can be an attractive proposition to get a Monte Carlo mechanism up and running quickly.

The process behind it is again very simple (Algorithm 6). Essentially, a sequence of dependent samples for a random vector $\mathbf{x} = \langle x_1, \dots, x_m \rangle$ are generated by sampling individually from the set of conditional distributions $P(x_i | \{x_j\}_{j \neq i})$; that is, we use the distributions of each element of \mathbf{x} conditioned on all other elements of \mathbf{x} . Furthermore, the conditional values used to generate the next state, \mathbf{x}_{i+1} , are supplied from the values of the current state, \mathbf{x}_i . This is done under that assumption that, individually, these distributions are easier to sample from than the joint distribution $P(\mathbf{x})$ as a whole.

A.6 Slice Sampling

Slice sampling (Mackay, 2003) is a MCMC approach that can be applied whenever the Metropolis-Hastings method can be applied, but it is more robust against choices in its parameters. The advantage is that it is self tuning: whereas Metropolis-Hastings is sensitive to the shape of its proposal density relative to the target, slice sampling adapts its parameters in response to the target distribution. This does not necessarily mean that increasing dimensionality does not decrease the efficiency of the algorithm, nor does the algorithm prevent random walking behaviour. However, it does mean that slice sampling will generally out perform basic Metropolis-Hastings techniques on a variety of problems.

Although the details of the algorithm are slightly more involved than those we have described so far, a full definition is not necessary to understand its main properties. In one dimension it can be considered similar to rejection sampling, in that it first generates 2 dimensional vectors sampled uniformly from under the curve of $P^*(x)$. However, no fixed proposal density is used. Instead, each subsequent sample $\langle x', u' \rangle$ is drawn based on $\langle x, u \rangle$ as illustrated in Figure A.5. First, $P^*(x)$ is evaluated, and then u' is drawn uniformly from the interval $[0, P^*(x)]$. Second, an interval $[x_l, x_u]$ that includes x is randomly selected, and from this x' is chosen. During this process, any proposed samples for which $u' > P^*(x')$ are rejected, and the interval $[x_l, x_u]$ is changed dynamically so that it does not extend too far beyond or within the probable region of $P(x)$. In this way, random walk behaviour is controlled, while at the same time the

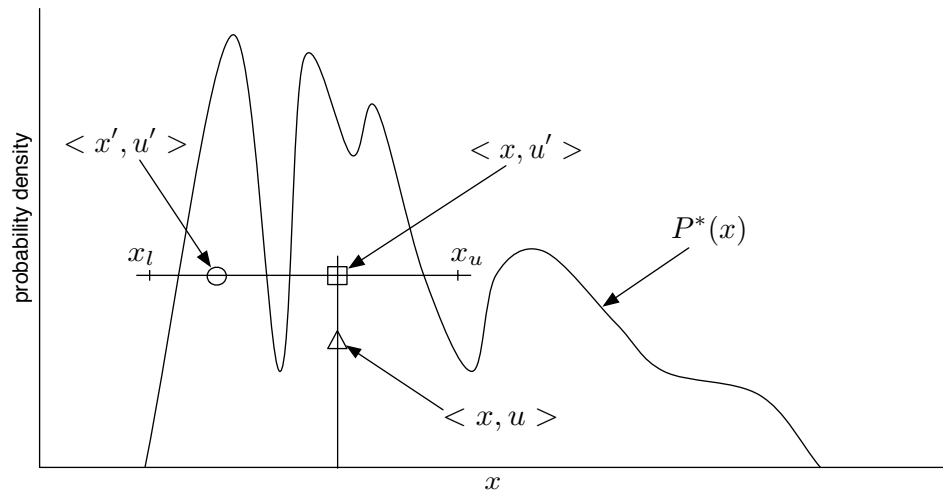


FIGURE A.5: Illustration of the slice sampling algorithm.

rejection rate is kept low, ensuring samples are unlikely to be proposed in areas of low density.

This basic algorithm for one-dimensional problems can also be modified to deal with multi-dimensional problems. However, perhaps more usefully, there is a computer-friendly version of the algorithm due to [Skilling and Mackay \(2003\)](#) that can be applied to problems of any number of dimensions without modification. To achieve this, Skilling and MacKay account for the fact that, no matter what the distribution's domain, real variables will always be represented using a finite number of bits. Therefore, all the operations in the algorithm are applied directly to a set of bits of given length, rather than to the real values that they conceptualise.

Appendix B

Parameter Mapping for the Beta Distribution

In this appendix, we provide two theorems, which show how the α and β parameters of a beta distribution can be calculated, if we know the variance and mean of the distribution. Specifically, Theorem B.1 shows how α can be derived in terms of the distribution mean (denoted μ) and the variance (denoted σ^2), and then how, given this, β can be determined from α and μ . Following this, Theorem B.2, gives an alternative expression for β , in terms of σ^2 and μ only.

Theorem B.1. *Given Equations B.1 & B.2, the parameters of the beta distribution, α & β can be derived from the distribution variance (denoted σ^2) and the mean (denoted μ).*

$$\mu = \frac{\alpha}{\alpha + \beta} \quad (\text{B.1})$$

$$\sigma^2 = \frac{\alpha \cdot \beta}{(\alpha + \beta)^2(\alpha + \beta + 1)} \quad (\text{B.2})$$

Proof: *First of all, we express β in terms of μ and α :*

$$\mu = \frac{\alpha}{\alpha + \beta} \quad (\text{from definition}) \quad (\text{B.3})$$

$$(\alpha + \beta) \cdot \mu = \alpha \quad (\text{B.4})$$

$$\alpha + \beta = \alpha / \mu \quad (\text{B.5})$$

$$\beta = \alpha / \mu - \alpha \quad (\text{B.6})$$

Now substitute for β in equation B.2 and simplify:

$$\sigma^2 = \frac{\alpha(\alpha/\mu - \alpha)}{(\alpha + (\alpha/\mu - \alpha))^2(\alpha + (\alpha/\mu - \alpha) + 1)} \quad (\text{B.7})$$

$$\sigma^2 = \frac{\alpha^2/\mu - \alpha^2}{(\alpha/\mu)^2(\alpha/\mu + 1)} \quad (\text{B.8})$$

$$\sigma^2 = \frac{\alpha^2/\mu - \alpha^2}{(\alpha/\mu)^3 + (\alpha/\mu)^2} \quad (\text{B.9})$$

$$\sigma^2 = \frac{\alpha^2/\mu - \alpha^2}{\alpha^3/\mu^3 + \alpha^2/\mu^2} \quad (\text{B.10})$$

$$\sigma^2 = \frac{1/\mu - 1}{\alpha/\mu^3 + 1/\mu^2} \quad (\text{B.11})$$

$$\sigma^2 = \frac{\mu^2 - \mu^3}{\alpha + \mu} \quad (\text{B.12})$$

Now arrange to find α :

$$\sigma^2(\alpha + \mu) = \mu^2 - \mu^3 \quad (\text{B.13})$$

$$\sigma^2 \cdot \alpha + \sigma^2 \cdot \mu = \mu^2 - \mu^3 \quad (\text{B.14})$$

$$\sigma^2 \cdot \alpha = \mu^2 - \mu^3 - \sigma^2 \cdot \mu \quad (\text{B.15})$$

$$\alpha = (\mu^2 - \mu^3 - \sigma^2 \cdot \mu) / \sigma^2 \quad (\text{B.16})$$

$$\alpha = \frac{\mu^2 - \mu^3}{\sigma^2} - \mu \quad (\text{B.17})$$

From Equations B.6 and B.17, α and β can be expressed as follows, thus proving the theorem.

$$\alpha = \frac{\mu^2 - \mu^3}{\sigma^2} - \mu, \quad \beta = \frac{\alpha}{\mu} - \alpha$$

Theorem B.2. *The β parameter of the beta distribution can be expressed only in terms of μ and σ^2 as shown in Equation B.18. We prove this in two ways: first, by considering the properties of the beta distribution; and second, by substitution.*

$$\beta = \frac{(1 - \mu)^2 - (1 - \mu)^3}{\sigma^2} - (1 - \mu) \quad (\text{B.18})$$

Proof through the properties of the Beta Distribution: *Imagine that we have two beta distributions: distribution d with parameters α and β , and distribution \hat{d} with parameters $\hat{\alpha}$ and $\hat{\beta}$. Similarly, we denote the mean of \hat{d} as $\hat{\mu}$ and the variance of \hat{d} as $\hat{\sigma}^2$.*

Now assume that $\hat{\alpha} = \beta$ and $\hat{\beta} = \alpha$. From this we know that $\hat{\sigma}^2 = \sigma^2$ since:

$$\frac{\alpha \cdot \beta}{(\alpha + \beta)^2(\alpha + \beta + 1)} = \frac{\beta \cdot \alpha}{(\beta + \alpha)^2(\beta + \alpha + 1)} = \frac{\hat{\alpha} \cdot \hat{\beta}}{(\hat{\alpha} + \hat{\beta})^2(\hat{\alpha} + \hat{\beta} + 1)} \quad (\text{B.19})$$

and $\hat{\mu} = (1 - \mu)$ since:

$$\hat{\mu} + \mu = \frac{\alpha}{\alpha + \beta} + \frac{\hat{\alpha}}{\hat{\alpha} + \hat{\beta}} \quad (\text{B.20})$$

$$\hat{\mu} + \mu = \frac{\alpha}{\alpha + \beta} + \frac{\beta}{\beta + \alpha} \quad (\text{B.21})$$

$$\hat{\mu} + \mu = 1 \quad (\text{B.22})$$

$$\hat{\mu} = 1 - \mu \quad (\text{B.23})$$

We can now prove Equation B.18 as follows:

$$\beta = \hat{\alpha} = \frac{\hat{\mu}^2 - \hat{\mu}^3}{\hat{\sigma}^2} - \hat{\mu}, \quad (\text{from Equation B.17}) \quad (\text{B.24})$$

$$\beta = \frac{(1 - \mu)^2 - (1 - \mu)^3}{\sigma^2} - (1 - \mu), \quad (\text{by substitution}) \quad (\text{B.25})$$

Proof by Substitution: We now show that Equation B.18 is true by substituting Equation B.17 into Equation B.6 as follows:

$$\beta = \frac{\alpha}{\mu} - \alpha \quad (\text{B.26})$$

$$\beta = \left[\frac{\mu^2 - \mu^3}{\sigma^2} - \mu \right] / \mu - \left[\frac{\mu^2 - \mu^3}{\sigma^2} - \mu \right] \quad (\text{B.27})$$

$$\beta = \left[\frac{\mu - \mu^2}{\sigma^2} - 1 \right] - \left[\frac{\mu^2 - \mu^3}{\sigma^2} - \mu \right] \quad (\text{B.28})$$

$$\beta = \frac{(\mu - \mu^2) - (\mu^2 - \mu^3)}{\sigma^2} - (1 - \mu) \quad (\text{B.29})$$

$$\beta = \frac{\mu - 2\mu^2 + \mu^3}{\sigma^2} - (1 - \mu) \quad (\text{B.30})$$

To show that Equations B.18 and B.30 are equivalent, we expand $(1 - \mu)^2 - (1 - \mu)^3$.

$$(1 - \mu)^2 = 1 - 2\mu + \mu^2 \quad (\text{B.31})$$

$$(1 - \mu)^3 = (1 - 2\mu + \mu^2)(1 - \mu) \quad (\text{B.32})$$

$$(1 - \mu)^3 = (1 - 2\mu + \mu^2) - (\mu - 2\mu^2 + \mu^3) \quad (\text{B.33})$$

$$(1 - \mu)^3 = 1 - 2\mu + \mu^2 - \mu + 2\mu^2 - \mu^3 \quad (\text{B.34})$$

$$(1 - \mu)^3 = 1 - 3\mu + 3\mu^2 - \mu^3 \quad (\text{B.35})$$

$$(1 - \mu)^2 - (1 - \mu)^3 = (1 - 2\mu + \mu^2) - (1 - 3\mu + 3\mu^2 - \mu^3) \quad (\text{B.36})$$

$$(1 - \mu)^2 - (1 - \mu)^3 = 1 - 2\mu + \mu^2 - 1 + 3\mu - 3\mu^2 + \mu^3 \quad (\text{B.37})$$

$$(1 - \mu)^2 - (1 - \mu)^3 = \mu - 2\mu^2 + \mu^3 \quad (\text{B.38})$$

$$\beta = \frac{(1 - \mu)^2 - (1 - \mu)^3}{\sigma^2} - (1 - \mu) = \frac{\mu - 2\mu^2 + \mu^3}{\sigma^2} - (1 - \mu) \quad (\text{B.39})$$

Hence Equations B.18 and B.30 are equivalent and therefore Equation B.18 is true.

Bibliography

- A. Abdul-Rahman and S. Hailes. **A distributed trust model**. In *Proceedings of the 1997 Workshop on New Security Paradigms*, pages 48–60, Cumbria, UK, September 1997. ACM Press.
- F. Adelstein, S. Gupta, R. Golden, and L. Schweibert. *Fundamentals of Mobile and Pervasive Computing*. McGraw-Hill, November 2004.
- R. Ashri, S. D. Ramchurn, J. Sabater, M. Luck, and N. R. Jennings. **Trust evaluation through relationship analysis**. In *Proceedings of the 4th International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 1005–1011, Utrecht, the Netherlands, July 2005. ACM Press.
- F. Azzedin and M. Maheswaran. **Evolving and managing trust in grid computing systems**. In *Proceedings of the IEEE Canadian Conference on Electrical and Computer Engineering*, volume 3, pages 1424–1429, Manitoba, Canada, May 2002a. IEEE Computer Society.
- F. Azzedin and M. Maheswaran. **Integrating trust into grid resource management systems**. In *Proceedings of the 2002 International Conference on Parallel Processing*, pages 47–54, British Columbia, Canada, August 2002b. IEEE Computer Society.
- F. Azzedin and M. Maheswaran. **Towards trust-aware management in grid computing systems**. In *Proceedings of the 2nd IEEE/ACM International Symposium on Cluster Computing and the Grid*, pages 419–424, Berlin, Germany, May 2002c. IEEE Computer Society.
- K. S. Barber and J. Kim. Belief revision process based on trust: Agents evaluating reputation of information sources. In Rino Falcone, Munindar P. Singh, and Yao-Hua Tan, editors, *Trust in Cyber-societies*, volume 2246 of *Lecture Notes in Computer Science*, pages 73–82. Springer-Verlag, 2001.
- T. Berners-Lee, J. Hendler, and O. Lassila. **The semantic web**. *Database and Network Journal*, 36(3):7–14, June 2006.

- B. Blankenburg, R. K. Dash, S. D. Ramchurn, M. Klusch, and N. R. Jennings. Trusted kernel-based coalition formation. In *Proceedings of the 4th International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 989–996, Utrecht, the Netherlands, July 2005. ACM Press.
- M. Blaze, J. Feigenbaum, and J. Lacy. **Decentralized trust management**. In *Proceedings of the 1996 IEEE Symposium on Security and Privacy*, pages 164–173, Oakland, California, May 1996. IEEE Computer Society.
- S. Buchegger and J. Y. Le Boudec. A robust reputation system for mobile ad-hoc networks ic/2003/50. Technical report, Laboratory for Computer Communications and Applications, Ecole Polytechnique Fédérale de Lausanne, 2003.
- R. L. Burden and J. D. Faires. *Numerical analysis*. Brookes/Cole Publishing, 8th edition, 2001.
- F. Campos and S. Cavalcante. **An extended approach for dempster-shafer theory**. In *Proceedings of the 2003 IEEE International Conference on Information Reuse and Integration*, pages 338–344, Las Vegas, NV, USA, October 2003. IEEE Computer Society.
- C. Castelfranchi and R. Falcone. Social trust: A cognitive approach. In C. Castelfranchi and Y. Tan, editors, *Trust and Deception in Virtual Societies*, pages 55–90. Kluwer Academic Publishers, 2001.
- P. R. Cohen. *Empirical Methods for Artificial Intelligence*. M.I.T. Press, 1995.
- E. Damsleth. Conjugate classes for gamma distributions. *Scandinavian Journal of Statistics*, 2(2):80–84, 1975.
- V. D. Dang and N. R. Jennings. Polynomial algorithms for clearing multi-unit single item and multi-unit combinatorial reverse auctions. In *Proceedings of the 15th European Conference on Artificial Intelligence*, pages 23–27, Lyon, France, July 2002. IOS Press.
- V. D. Dang and N. R. Jennings. Coalition structure generation in task-based settings. In *Proceedings of the 17th European Conference on Artificial Intelligence*, pages 567–571, Trento, Italy, August 2006. IOS Press.
- P. Dasgupta. **Trust as a commodity**. In D. Gambetta, editor, *Trust: Making and Breaking Cooperative Relations*, chapter 4, pages 49–72. Basil Blackwell, 1988. Reprinted in electronic edition from Department of Sociology, University of Oxford, 2000.
- R. K. Dash, D. C. Parkes, and N. R. Jennings. **Computational mechanism design: A call to arms**. *IEEE Intelligent Systems*, 18(6):40–47, 2003.
- R. K. Dash, S. D. Ramchurn, and N. R. Jennings. **Trust-based mechanism design**. In *Proceedings of the 3rd International Joint Conference on Autonomous Agents and*

- Multiagent Systems*, pages 748–755, New York, USA, July 2004. IEEE Computer Society.
- M. DeGroot and M. Schervish. *Probability & Statistics*. Addison-Wesley, 3rd edition, 2002.
- C. Dellarocas. **Mechanisms for coping with unfair ratings and discriminatory behavior in online reputation reporting systems**. In *Proceedings of the 21st International Conference on Information Systems*, pages 520–525, Brisbane, Australia, December 2000.
- D. G. T. Denison, C. C. Holmes, B. K. Mallick, and A. F. M. Smith. *Bayesian Methods for Nonlinear Classification and Regression*. Wiley, 2002.
- V. Deora, J. Shao, W. A. Gray, and N. J. Fiddian. A quality of service management framework based on user expectations. In *Proceedings of the 1st International Conference on Service Oriented Computing*, volume 2910 of *Lecture Notes in Computer Science*, pages 104–114, Trento, Italy, December 2003. Springer-Verlag.
- V. Deora, J. Shao, G. Shercliff, P.J. Stockreisser, W.A. Gray, and N.J. Fiddian. Incorporating QoS specifications in service discovery. In *Proceedings of 2nd International Web Services Quality Workshop*, pages 252–263, Brisbane, Australia, November 2004. Springer-Verlag.
- Z. Despotovic and K. Aberer. **Probabilistic prediction of peers' performance in p2p networks**. *Engineering Applications of Artificial Intelligence*, 18(7):771–780, October 2005.
- P. S. Dutta, S. Dasmahapatra, S. R. Gunn, N. R. Jennings, and L. Moreau. Cooperative information sharing to improve distributed learning. In *Proceedings of the 3rd International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 828–835, New York, USA, July 2004. IEEE Computer Society.
- A. Erdélyi, W. Magnus, F. Oberhettinger, and F. G. Tricomi. *Higher-Transcendental Functions*, volume 1. McGraw-Hill, 1953.
- M. Evans, N. Hastings, and B. Peacock. *Statistical Distributions*. Wiley, 3rd edition, 2000.
- M. Evans and T. Swartz. *Approximating Integrals via Monte Carlo and Deterministic Methods*. Oxford University Press, 1999.
- R. Falcone, G. Pezzulo, and C. Castelfranchi. **Fuzzy approach to a belief-based trust computation**. In R. Falcone, S. Barber, L. Korba, and M. Singh, editors, *Trust, Reputation and Security: Theories and Practice*, volume 2631 of *Lecture Notes in Artificial Intelligence*, pages 73–86. Springer, 2003.

- M. Fan, Y. Tan, and A. B. Whinston. **Evaluation and design of online cooperative feedback mechanisms for reputation management.** *IEEE Transactions on Knowledge and Data Engineering*, 17(2):244–254, February 2005.
- I. Foster, N. R. Jennings, and C. Kesselman. Brain meets brawn: Why grid and agents need each other. In *Proceedings of the 3rd International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 8–15, New York, USA, July 2004. IEEE Computer Society.
- I. Foster and C. Kesselman, editors. *The Grid: Blueprint for a New Computing Infrastructure*. Morgan Kaufmann, 2nd edition, 2004.
- I. Foster, C. Kesselman, and S. Tuecke. The anatomy of the grid: Enabling scalable virtual organizations. *International Journal of Supercomputer Applications*, 15(3): 200–222, 2001.
- D. Gambetta. **Can we trust trust?** In D. Gambetta, editor, *Trust: Making and Breaking Cooperative Relations*, chapter 13, pages 213–237. Basil Blackwell, 1988. Reprinted in electronic edition from Department of Sociology, University of Oxford.
- M. Garrido. *Modélisation des évènements rares et estimation des quantiles extrêmes, Méthodes de sélection de modèles pour les queues de distribution*. PhD thesis, Université Joseph Fourier de Grenoble, France, 2002.
- J. E. Gentle. *Random number generation and Monte Carlo methods*. Springer-Verlag, 1998.
- E. H. Gerding, R. K. Dash, D. C. K. Yuen, and N. R. Jennings. Optimal bidding strategies for simultaneous vickrey auctions with perfect substitutes. In *Proceedings of the 8th International Workshop on Game Theoretic and Decision Theoretic Agents*, pages 10–17, Hakodate, Japan, 2006.
- D. Gollmann. *Computer Security*. Wiley, 1998.
- T. Grandison and M. Sloman. **A survey of trust in internet applications.** *IEEE Communications Surveys and Tutorials*, 3(4):2–16, 2000.
- N. Griffiths and K. Chao. **Experience-based trust: Enabling effective resource selection in a grid environment.** In P. Hermann, V. Issarny, and S. Shiu, editors, *Proceedings of the 3rd International Conference on Trust Management*, volume 3477 of *Lecture Notes in Computer Science*, pages 240–255, Rocquencourt, France, May 2005. Springer-Verlag.
- M. He, A. Rogers, X. Luo, and N. R. Jennings. Designing a successful trading agent for supply chain management. In *Proceedings of the 5th International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 1159–1166, Hakodate, Japan, 2006. ACM Press.

- M. Huhns, M. Singh, M. Burstein, K. Decker, E. Durfee, T. Finin, L. Gasser, H. Goradia, N. R. Jennings, K. Lakartaju, H. Nakashima, V. Parunak, J. Rosenschein, A. Ruvinsky, G. Sukthankar, S. Swarup, K. Sycara, M. Tambe, T. Wagner, and L. Zavala. Research directions for service-oriented multi-agent systems. *IEEE Internet Computing*, 9(6):65–70, 2005.
- M. N. Huhns and M. P. Singh. **Service-oriented computing: key concepts and principles**. *IEEE Internet Computing*, 9(1):75–81, 2005.
- A. Jøsang. **A logic for uncertain probabilities**. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 9(3):279–311, June 2001.
- A. Jøsang. **Subjective evidential reasoning**. In *Proceedings of the 9th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems*, pages 1671–1678, Annecy, France, July 2002.
- A. Jøsang and R. Ismail. **The beta reputation system**. In *Proceedings of the 15th Bled Conference on Electronic Commerce*, pages 324–337, Bled, Slovenia, June 2002.
- A. Jøsang, R. Ismail, and C. Boyd. **A survey of trust and reputation systems for online service provision**. *Decision Support Systems*, 2005. (In Press) Retrieved by <http://dx.doi.org/10.1016/j.dss.2005.05.019>.
- R. Jurca and B. Faltings. Towards incentive-compatible reputation management. In R. Falcone, S. Barber, L. Korba, and M. Singh, editors, *Trust, Reputation and Security: Theories and Practice*, volume 2631 of *Lecture Notes in Artificial Intelligence*, pages 138–147. Springer-Verlag, 2003.
- R. L. Keeney and H. Raiffa. *Decisions with Multiple Objectives: Preferences and Value Tradeoffs*. Wiley, 1976.
- A. I. Khuri. *Advanced Calculus with Applications in Statistics*. Wiley, 2nd edition, 2003.
- T. Klos and H. La Poutré. Using reputation-based trust for assessing agent reliability. In *Proceedings of the 7th International Workshop on Trust in Agent Societies*, pages 75–82, New York, USA, 2004.
- N. Littlestone and M. K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108(2):212–261, 1994.
- D. J. C. Mackay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003.
- S. Marsh. *Formalising Trust as a Computational Concept*. PhD thesis, University of Stirling, 1994.
- S. A. McIlraith, T. C. Son, and H. Zeng. **Semantic web services**. *IEEE Intelligent Systems*, 16(2):46–53, 2001.

- B. Misztal. *Trust in Modern Societies: The Search for the Bases of Social Order*. Polity Press, 1996.
- L. Mui, M. Mohtashemi, C. Ang, P. Szolovits, and A. Halberstadt. **Ratings in distributed systems: A bayesian approach**. In *Proceedings of the 11th Workshop on Information Technologies and Systems*, New Orleans, USA, December 2001.
- D. Nguyen, S. G. Thompson, J. Patel, W. T. L. Teacy, N. R. Jennings, M. Luck, V. Dang, S. Chalmers, N. Oren, T. J. Norman, A. Preece, P. M. D. Gray, G. Shercliff, P. J. Stockreisser, J. Shao, W. A. Gray, and N. J. Fiddian. Delivering services by building and running virtual organisations. *BT Technology Journal*, 24(1):141–152, 2006.
- T. D. Nguyen and N. R. Jennings. Managing commitments in multiple concurrent negotiations. *International Journal of Electronic Commerce Research and Applications*, 4:362–376, 2005.
- T. J. Norman, A. Preece, S. Chalmers, N. R. Jennings, M. Luck, V.D. Dang, T. D. Nguyen, V. Deora, , J. Shao, A. Gray, and N. J. Fiddian. Conoise: Agent-based formation of virtual organisations. In *Proceedings of 23rd SGAI International Conference on Innovative Techniques and Applications of AI*, pages 353–366, Cambridge, UK, 2003.
- R. H. J. M. Otten and L. P. P. P. van Ginneken. *The Annealing Algorithm*. Kluwer Academic Publishers, 1989.
- D. Pardoe, P. Stone, and M. Van Middlesworth. Tactex-05: An adaptive agent for tac scm. In *In Proceedings of the AAMAS-06 Workshop on Trading Agent Design and Analysis*, Hakodate, Japan, 2006.
- J. Patel, W. T. L. Teacy, N. R. Jennings, and M. Luck. A probabilistic trust model for handling inaccurate reputation sources. In *Proceedings of the 2nd European Workshop on Multiagent Systems*, pages 521–529, Barcelona, Spain, December 2004.
- J. Patel, W. T. L. Teacy, N. R. Jennings, and M. Luck. **A probabilistic trust model for handling inaccurate reputation sources**. In P. Hermann, V. Issarny, and S. Shiu, editors, *Proceedings of the 3rd International Conference on Trust Management*, volume 3477 of *Lecture Notes in Computer Science*, pages 193–209, Rocquencourt, France, May 2005a. Springer-Verlag.
- J. Patel, W. T. L. Teacy, N. R. Jennings, M. Luck, S. Chalmers, N. Oren, T. J. Norman, A. Preece, P. M. D. Gray, G. Shercliff, P. J. Stockreisser, J. Shao, W. A. Gray, N. J. Fiddian, and S. Thompson. **Agent-based virtual organisations for the grid**. In *Proceedings 1st International Workshop on Smart Grid Technologies*, pages 1–15, Utrecht, Netherlands, July 2005b.

- J. Patel, W. T. L. Teacy, N. R. Jennings, M. Luck, S. Chalmers, N. Oren, T. J. Norman, A. Preece, P. M. D. Gray, G. Shercliff, P. J. Stockreisser, J. Shao, W. A. Gray, N. J. Fiddian, and S. Thompson. Monitoring, policing and trust for grid-based virtual organisations. In *Proceedings of the 2005 UK OST e-Science All Hands Meeting*, pages 891–898, Nottingham, UK., September 2005c. EPSRC.
- J. Patel, W. T. L. Teacy, N. R. Jennings, M. Luck, S. Chalmers, N. Oren, T. J. Norman, A. Preece, P. M. D. Gray, G. Shercliff, P. J. Stockreisser, J. Shao, W. A. Gray, N. J. Fiddian, and S. Thompson. **Agent-based virtual organisations for the grid.** *International Journal of Multiagent and Grid Systems*, 1(4):237–249, 2006.
- C. P. Pfleeger. *Security in Computing*. Prentice Hall, 2002.
- S. Ramchurn, C. Sierra, L. Godo, and N. R. Jennings. A computational trust model for multi-agent interactions based on confidence and reputation. In *Proceedings of the 6th International Workshop of Deception, Fraud and Trust in Agent Societies*, pages 69–75, Melbourne, Australia, July 2003. ACM Press.
- J. Roy and A. Ramanujan. **Understanding web services.** *IT Professional*, 3(6):69–73, 2001.
- S. Russell and P. Norvig. *Artificial Intelligence A Modern Approach*. Prentice Hall, 2nd edition, 2003.
- J. Sabater and C. Sierra. **Regret: A reputation model for gregarious societies.** In *Proceedings of the 5th International Conference on Autonomous Agents*, pages 194–195, Montreal, Canada, 2001.
- J. Sabater and C. Sierra. Social regret, a reputation model based on social relations. *SIGecom Exchanges*, 3(1):44–56, 2002.
- G. Shafer. *A Mathematical Theory of Evidence*. Princeton University Press, 1976.
- J. Shao, W. A. Gray, N. J. Fiddian, V. Deora, G. Shercliff, P. J. Stockreisser, T. J. Norman, A. Preece, P. M. D. Gray, S. Chalmers, N. Oren, N. R. Jennings, M. Luck, V. D. Dang, T. D. Nguyen, J. Patel, and W. T. L. Teacy. **Supporting formation and operation of virtual organisations in a grid environment.** In *Proceedings of the 2004 UK OST e-Science All Hands Meeting*, pages 376–383, Nottingham, UK., September 2004. EPSRC.
- C. Sierra. Agent-mediated electronic commerce. *Autonomous Agents and Multi-Agent Systems*, 9(3):285–301, 2004.
- J. Skilling and D. J. C. Mackay. **Slice sampling — a binary implementation.** *Annals of Statistics*, 31(3):753–755, 2003. Discussion of *Slice Sampling* by R. M. Neal.

- L. Sun, L. Jiao, Y. Wang, S. Cheng, and W. Wang. **An adaptive group-based reputation system in peer-to-peer networks**. In *Proceedings of the 1st International Workshop on Internet and Network Economics*, volume 3828 of *Lecture Notes in Computer Science*, pages 651–659, Hong Kong, China, December 2005. Springer-Verlag.
- R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, 1998.
- W. T. L. Teacy, J. Patel, N. R. Jennings, and M. Luck. **Coping with inaccurate reputation sources: Experimental analysis of a probabilistic trust model**. In *Proceedings of 4th International Joint Conference on Autonomous Agents and MultiAgent Systems*, pages 997–1004, Utrecht, the Netherlands, July 2005. ACM Press.
- W. T. L. Teacy, J. Patel, N. R. Jennings, and M. Luck. **Travos: Trust and reputation in the context of inaccurate information sources**. *Autonomous Agents and Multi-Agent Systems*, 12(2):183–198, February 2006.
- T. Tran and R. Cohen. Improving user satisfaction in agent-based electronic marketplaces by reputation modelling and adjustable product quality. In *Proceedings of the 3rd International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 828–835, New York, USA, July 2004. IEEE Computer Society.
- G. Upton and I. Cook. *Dictionary of Statistics*. Oxford University Press, 2002.
- Y. Wang and J. Vassileva. **Bayesian network-based trust model**. In *Proceedings of IEEE/WIC International Conference on Web Intelligence*, pages 372–378, Halifax, Canada, October 2003.
- A. Whitby, A. Jøsang, and J. Indulska. **Filtering out unfair ratings in bayesian reputation systems**. In *Proceedings of the 7th International Workshop on Trust in Agent Societies*, pages 106–117, New York, USA, 2004.
- M. J. Wooldridge and N. R. Jennings. **Software engineering with agents: pitfalls and pratfalls**. *IEEE Internet Computing*, 3(3):20–27, 1999.
- B. Yu and M. P. Singh. An evidential model of distributed reputation management. In *Proceedings of 1st International Joint Conference on Autonomous Agents and Multi-agent Systems*, volume 1, pages 294–301, Bologna, Italy, July 2002. ACM Press.
- Bin Yu and Munindar P. Singh. **Detecting deception in reputation management**. In *Proceedings of the 2nd International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 73–80, Melbourne, Australia, July 2003. ACM Press.
- G. Zacharia, A. Moukas, and P. Maes. Collaborative reputation mechanisms in electronic marketplaces. In *Proceedings of 32nd Hawaii International Conference on System Sciences*, Maui, Hawaii, January 1999. IEEE Computer Society Press.

-
- L. A. Zadeh. Fuzzy sets. *Information and Control*, 8(3):338–353, June 1965.
- L. A. Zadeh. Fuzzy logic and approximate reasoning. *Synthese*, 30:407–428, 1975.
- J. Zhang and R. Cohen. A personalized approach to address unfair ratings in multiagent reputation systems. In *Proceedings of the 9th International Workshop on Trust in Agent Societies*, pages 89–98, May 2006.