



Data for monitoring the SDGs need to be at fine spatial and temporal scales to enable decision makers and researchers to track and understand the trajectories in development progress (1). Mismatches in scale could be a problem for understanding socio-ecological systems (18) because human uses of, and dependencies on, natural resources may differ depending on the scale at which analysis is performed (19). Past studies have highlighted the potential for RS data to be used for poverty mapping at aggregated community levels such as the village (9), groups of villages (8), or census enumeration districts (7). Aggregating household and landscape information can result in the modifiable areal unit problem (20), due to the need to construct artificial boundaries. This effectively means that the same set of data can produce different results depending on how data are aggregated and lead to erroneous conclusions. In general, the average values from single polygons used to link RS and socioeconomic data in the past mask the multilevel interactions that occur between households and environmental resources. Aggregating environmental resources into a single polygon covering multiple households assumes that all households have the same opportunity to use the landscape to pursue livelihood strategies. This could have substantial consequences for policy recommendations based on understandings of the relationship between wealth and environment resulting from these analyses (21). Wealth can vary between neighboring households. Therefore, it is reasonable to expect that the relationships between wealth and RS features will differ at the community and household level. To examine these complex relationships requires analysis of wealth and RS features at finer spatial scales than done previously.

Fine spatial resolution satellite data could be helpful for monitoring SDG1 “Ending Poverty”; in particular, it could contribute to identifying extreme poverty and those areas likely affected by poverty, targeting resource allocation, and building rural resilience to climatic and environmental impacts. In this study, we hypothesized that, fine-grained socio-economic and environmental data allow a more mechanistic understanding of human–environment interactions. We tested this hypothesis using a case study in rural Kenya by predicting household level wealth using environmental characteristics extracted from RS data. We examine two study questions crucial to understand whether RS data can be used to bridge the data gaps in monitoring aspects of household wealth: (i) Can the variance in household wealth be explained with RS data? (ii) Does a socioecologically informed approach to treating RS data increase the ability to explain the variance in household level wealth?

## Results

We used a classification tree to examine if RS data could be used to predict household level wealth in the rural village of Sauri, Kenya. Within the study area, households typically live in homesteads, small areas with several structures, gardens or woodlots, and a surrounding hedge. Agricultural fields are interspersed between homesteads. Agriculture is the primary livelihood, with maize the main crop and bananas, beans, cassava, kale, and sorghum also grown. Rainfall is bimodal, allowing two cropping seasons: the long rains (March–June) during which the majority of maize crops are grown and the short rains (September–December), which are highly variable. This area is typical of many small-holder farming landscapes in East Africa; it is highly fragmented, densely populated, and topographically varied, with a complex mosaic of land cover classes. In 2005, 79% of the Sauri population was living below \$1 per day (1993 PPP) and 89.5% below \$2 per day (22).

We developed a multilevel approach to examine the relationships between household wealth and RS features at four spatial levels: level 1 homestead, level 2 agricultural land, level 3 village cluster, and level 4 wider village periphery (Fig. 1 and described in *SI Appendix*, Fig. S1). This method was compared with the single-level approach previously used for predicting wealth with aggregated socioeconomic data. Overall model accuracy for the multilevel approach was 60% using the training data and 45% using the testing data, between 6 and 12% higher than that using the single-level approach (Table 1). The

predictive accuracy for explaining the variance in the poorest households increased from 52% in the single-level approach to 62% using the multilevel approach. *t* tests indicated that the overall test accuracy and accuracy of wealth group 1 were significantly different between multilevel and single-level approaches (*SI Appendix*, Table S3).

The statistical relationships between household level wealth and multilevel RS features are shown in Fig. 2. The most important predictor variable appears at the top of the tree, meaning that building size was the most important RS variable for explaining the variance in household wealth. Other important variables in decreasing order of importance were amount of bare agricultural land and planted agricultural land adjacent to the homestead (level 2), amount of bare land in the homestead (level 1), the count of years that the number of agricultural growing days was lower than the 14-y average for that pixel, the growing period for year 2005 of the HH survey (level 4) and the amount of land classed as homestead within the common pool resource buffer (level 3).

The poorest households were characterized by a small building size (level 1), a relatively large proportion (almost half) of bare agricultural land in level 2 and bare ground in level 1 (Fig. 2). If a household had less than 43% bare ground within the homestead area, but with less than 163 growing days in the year, it was classified in the poorest household category. Poor households that had a large building size (37/92) had less than 21% of the agricultural land planted in September, but experienced over 6 y of below-average growing periods during the 14-y time series of Normalised Difference Vegetation Index (NDVI) and had over 16% of the common pool resource buffer (level 3) covered in homestead areas. Overall, 60% (55 households of a total of 92) of group 1 households, 31% (29 households of 92) of group 2 households and only 9% of group 3 households had a building size under 140 m<sup>2</sup>.



**Fig. 1.** The multilevel approach to linking households and landscape characteristics. Households have individual access to homestead areas (A, B, and C: level 1) and agricultural fields (A1–A3; B1–B3, C1–C3: level 2) surrounding the homestead. These levels should be linked to a single household. Households will also make use of common pool resources (level 3) around the village, which can be linked to multiple households. The wider regional level (level 4) considers infrastructure access. X, Y, and Z indicate fields that are adjacent to multiple households or no households, which would be split using our current method.

Approach	Tree size	Test accuracy, %	Training accuracy, %	Group 1, %	Group 2, %	Group 3, %
Multilevel	7.7	45	59	62	51	55
Single-level	10.4	38	59	50	49	52

Results are averaged from 1,000 iterations of the model trained on 80% of the household sample and tested using the remaining 20%. Group 1 is the poorest 40% of households, group 2 the middle 40%, and group 3 the wealthiest 20% of households.

The majority of wealthy households were characterized as having a large building size ( $>140\text{ m}^2$ ), less than 21% of the agricultural area planted by September 2004—the beginning of the short rainy season, more than 6 y of below average growing period, and less than 16% of the level 3 common pool resource area classed as homestead. Wealthy households with a small building size only had a small amount of unplanted agricultural land within the agricultural fields (level 3).

## Discussion

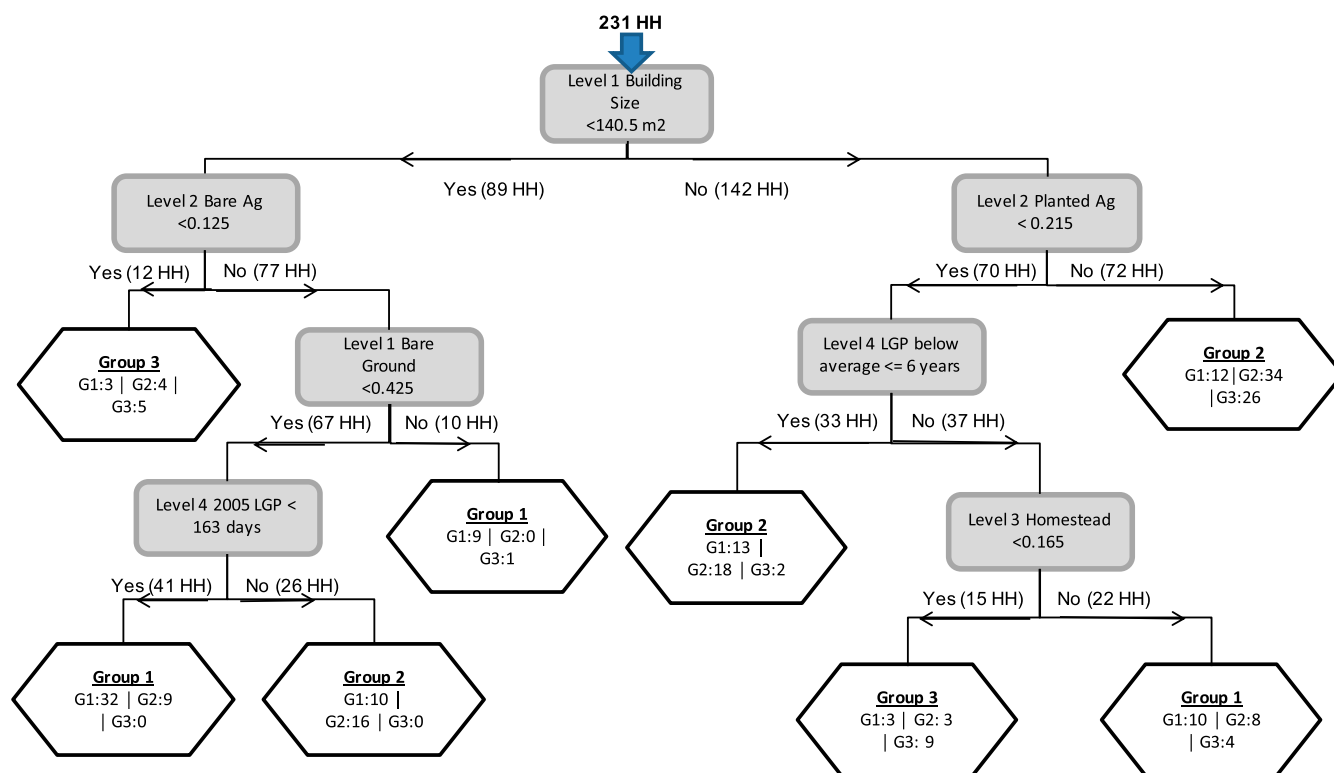
The multilevel approach included more complex types of land use and resource access based on the spatial arrangement of homesteads and agricultural fields, compared with a traditional single-level analysis. Our results show that considering socioecological conditions at multiple levels increases the accuracy of predicting wealth from RS data.

**Can Household Wealth Be Predicted from RS Data?** This study considers if wealth can be predicted from RS data at the household level. Predicting wealth in this area from RS data using a multilevel approach had an overall accuracy of 45% averaged over

1,000 model iterations. This is similar to past studies that predicted socioeconomic outcomes from RS data at coarser spatial resolutions (7–9). However, the multilevel approach developed here explained 62% of the variation in household wealth for the poorest group. A relatively high accuracy considering the complexities of household wealth and predictor variables that were derived from a single satellite image.

### Does a Multilevel Treatment Increase the Ability To Explain Variances in Household Poverty? The multilevel approach maps homestead

**in Household Poverty?** The multilevel approach maps homestead characteristics, local land uses, and agricultural productivity and relates them to a single household. Results indicate that splitting the RS features into different levels can have a positive impact on model accuracy as the optimal classification trees used features derived from all four levels (Fig. 2). There was a 10% increase in predictive capacity between the multilevel and single-level approaches for group 1, but little or no difference when predicting groups 2 and 3. Wealthier households may be less reliant on agriculture for food and income with nonfarm incomes such as salaries, business enterprises, and remittances contributing more to income in wealthier Kenyan households.



**Fig. 2.** Tree derived from cross-validation with an overall classification accuracy of 52%. Brackets after Yes/No indicate the number of households (HH) that met the split criteria. Group 1 = poorest, group 2 = middle, and group 3 = wealthiest households correspond to the predicted wealth group using the preceding data splits. G1/G2/G3 indicate the number of households observed in each wealth group at that terminal node. LGP, length of growing season. Level 1, homestead; level 2, agriculture; level 3, common-pool resource area; level 4, wider region for accessibility and length of growing period; bare ag, proportion of bare agricultural land within level 2.



The single-level approach assumes that all land within the buffer zone can be accessed and utilized by a given household. If an RS feature appears in multiple buffer zones, it will be linked to multiple households (Fig. 3), while in reality access to resources may be restricted to a single household. For example, homestead areas will most likely only be used by the household embedded within it. Of the 1,150 homesteads in the study area, 1,149 had more than one overlapping buffer zone with an average of 17 overlaps and maximum of 38. Thus, RS features within a homestead, which should only be linked to a single household, could be associated with up to 37 different households when using the single-level approach. This risks misestimating many households' resource access and introduces error into predictive models. The multiscale method can account for common pool resources such as hedges that are accessed by multiple households and separate them from agricultural fields and homesteads, which are likely used by single households. This result indicates that work using open data with displaced GPS coordinates such as that available from the Demographic and Health Surveys (DHS) may not be as useful for monitoring socioecological systems at fine spatial resolutions.

**Relationships between RS variables and household wealth.** The most important variables for explaining variance in household wealth were size of the household's buildings (level 1) and proportion of agriculture and bare land in level 2 (Fig. 2). The majority of households with small building sizes were from the poorest wealth categories (*SI Appendix, Table S2*). Small buildings likely indicate that a household has limited financial capital stock or has a small family size (human capital) with reduced labor pool and a lower diversity of livelihood strategies. Building size is not a seasonally dependent variable and could therefore provide a consistent RS variable for predicting rural wealth. The small number of households that had a small building size and were from the wealthiest group were differentiated from the poorer households by having a relatively small amount of bare agricultural land surrounding the homestead (level 2 nonvegetated <12.5%).

Tree regression allows for complexities to be identified in the relationships between wealth and RS variables. Households characterized by large building sizes had a lower proportion of bare agricultural land and a lower proportion of planted agricultural land at the start of the second planting season. Wealthier households derived 71% of their incomes from nonagricultural sources (23). Therefore, the results may indicate that these households do not need to plant second crops during the short rainy season. Poorer households were characterized as having more bare land within the agricultural fields (level 2) in September, which means

that the land has likely been prepared for planting for the short rains. This is an important finding because planting during the short rains is a high-risk strategy as around 50% of harvests fail due to drought (23). This result is consistent with poorer households planting second crops through necessity due to a lack of options for growing food or generating incomes (21% of the poorest households income was derived from nonfarm activities) (23).

The main growing period in the study area is around 155 d long (between March and July) and Moderate Resolution Imaging Spectro-radiometer (MODIS) data indicate a double cropping pattern. Therefore, the model prediction that poorer households had a total growing period of <163 d is indicative of two short agricultural seasons. This could be because poorer households delay planting while hiring themselves out to plant other farmers' fields for cash payments that are used to fund their own planting. This would result in late planting and a shorter main growing period compared with wealthier households. However, it could also be due to poorer households planting different crops with different maturing periods.

A large proportion of bare ground within a homestead (level 1) was associated with the poorest households. While it cannot be determined from the imagery, bare ground in the homestead would have different uses in different homesteads. Households use this space for socializing, and drying crops among other uses. Field observations indicated that wealthier households were more likely to invest in "greening" the homesteads to provide fencing poles, wind breaks, and pasture.

**The role of remote sensing in the data revolution for SDG monitoring programs.** The increasing availability of high-resolution satellite data means that methods, such as those developed in this study, could support the SDG "data revolution" (4) and provide a more cost-effective way of monitoring development than annual household surveys. The World Bank estimates the costs for a household survey at \$322.99 (USD 2014 prices) per household in Sub-Saharan Africa (24). This is the gold standard for surveys as it includes multiple modules and household visits. If the World Bank cost estimates were used to collect the socioeconomic information of the 330 households originally surveyed in our study site in Sauri, the total cost would be in the region of \$106,500 per year. In comparison, acquisition of high-resolution satellite imagery for the 100-km<sup>2</sup> site ranged from \$1,750 to \$5,000 per year (*SI Appendix, Table S4*). The World Bank proposes to survey countries every 3 y using sample surveys of between 3,000 and 10,000 households depending on the country (24). However, to monitor socioeconomic conditions sufficiently, some form of annual survey is recommended (3). Therefore, the World Bank approach leaves up to 10 y during the 15-y SDG period with no household surveys during the SDG timeframe, which could risk our understanding of the dynamics of change. If the sampled households are a panel, satellite data covering these households could be acquired every year to provide continual monitoring of some socioecological conditions and potentially provide \$100,000s worth of savings compared with household survey costs.

**Future Work.** The methodology developed here would need to be tested in multiple places, with different spatial arrangements of homes and agricultural fields, configurations of common resource areas, road networks, and market access. The approach still lacks detailed land tenure information but could vary the size of level 3 based on land ownership. Not all households have the same access to land and common pool resources across the landscape (25). Local and regional institutions can also impact the ways different actors access and utilize natural resources (26). Therefore, future work should examine how protected land areas, tenure rights, and institutional arrangements could be integrated into the multilevel approach. This could result in more accurate links between individual households and the parcels of land which they use. Developments in data and technology availability since the household survey was collected in 2005 provide significant future opportunities that should be explored for mapping wealth, health, and life on land. We



**Fig. 3.** The single-level approach to linking satellite and household data often uses a single radial buffer zone. This can be problematic as it results in overlapping regions and multiple pixels being assigned to multiple households when households would not have access to some land parcels such as multiple homesteads.





from January 2001 to December 2006. Each 200-m buffer zone was linked to the 500-m MODIS pixel in which it was contained; if a buffer zone was on the boundary of two or more pixels, it was given the average value. A Savitzky-Golay filter with a window size of six was used to smooth the data to estimate the length of growing period per year for each pixel (*SI Appendix, section S3*). The growing period was defined as the sum of the length of both growing periods in each year. Season start and end points were identified as the point where NDVI increased/decreased by 10% of the distance between the minimum and maximum and was computed in the TIMESAT software (35).

**Multilevel Approach.** We developed a mechanistic approach to represent the complexity of land and resource availability by considering capital endowments used exclusively by single households, those used by multiple households, common pool resources, and community infrastructure (*SI Appendix, section S1* for more details). For this case, we identified four levels highlighted in a stylized landscape model in Fig. 1 and *SI Appendix, Fig. S2*. At each level, particular RS features are extracted, e.g., land use within the area, vegetation productivity, and access measures such as distance to roads and market. Unless otherwise stated, land use proportions were extracted at each level using the isectpolyst tool in the GME.

**Predicting Household Wealth Using Remotely Sensed Features.** Household wealth was predicted with RS predictor variables using classification trees in R 3.3.2 (R Development Core Team 2016) and the “tree” package (36). (Some analytical steps could only be performed in a single software package at the time of analysis and so multiple software packages were used. eCognition, the only software allowing for multilevel object-based image classification and region growing of the homesteads; ESRI ArcMap, industry standard GIS software; GME tool, was able to deal with overlapping radial buffers which is not possible in ArcMap Buffer tool.) Classification trees have several benefits for this type of analysis. They are simple to implement and interpret and do not assume a normal error distribution. Classification trees are also hierarchical, allowing each variable to be used for splits multiple times (37),

effectively meaning that nonlinear relationships can be handled, important for modeling population–environment relationships (9). To reduce the problem of overfitting, we split the data into training/calibration (80% of the total data) and testing/validation (20% of the total data) samples (37). Each of the three wealth groups were sampled independently to ensure that the testing dataset contained 40% of households from the poorest wealth group, 40% from the middle, and 20% from the wealthiest group. The optimal tree was identified using a cross-validation approach, which prevents the model algorithm overfitting and predicting random noise in the data. The full tree was pruned to the size of the optimal tree; pruning is an essential step for generating useful predictions and ensures the most parsimonious tree with the highest predictive accuracy is obtained. The  $y$  variable was the wealth group of the household (1–3), and the  $x$  variables were the various RS features (*SI Appendix, Table S1*). The model was applied to the testing sample and a confusion matrix created using the “caret” package (38) to identify the overall model prediction accuracy as well as the accuracy of each wealth group. We repeated this process 1,000 times with the seed changed in each iteration to ensure a different set of households were included in the training and testing samples. This number of iterations ensured convergence in the calculated model prediction accuracies. This process was repeated for models using RS features extracted from the single-level approach and the multilevel approach for comparison.

**ACKNOWLEDGMENTS.** We thank Dr. Mark Musumba and Prof. Mat Williams for comments on an earlier version of the manuscript and Sombras Blancas Art and Design for converting Fig. 1 to digital graphics. The project received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie Grant agreement 656811. J.-C.S. considers this work a contribution to the Danish National Research Foundation Niels Bohr professorship project Aarhus University Research on the Anthropocene and a contribution to his VILLUM Investigator project “Biodiversity Dynamics in a Changing World” (Grant 16549). C.A.P. considers this work a contribution to the Millennium Villages Project, from which all the household data and high-resolution images were obtained.

- Griggs D, et al. (2013) Policy: Sustainable development goals for people and planet. *Nature* 495:305–307.
- Jacob A (2017) Mind the gap: Analyzing the impact of data gap in Millennium Development Goals' (MDGs) indicators on the progress toward MDGs. *World Dev* 93: 260–278.
- Jerven M (2014a) Benefits and costs of the data for development targets for the Post-2015 Development Agenda, Data for Development Assessment Working Paper, September 16, 2014.
- IEAG (2014) A world that counts: Mobilising the data revolution for sustainable development. *Independent Expert Advisory Group on a Data Revolution for Sustainable Development* (United Nations, New York).
- Devarajan S (2013) Africa's statistical tragedy. *Rev Income Wealth* 59:9–15.
- Jerven M (2014b) Poor numbers and what to do about them. *Lancet* 383:594–595.
- Engstrom R, Hersh J, Newhouse D (2016) Poverty in HD: What does high resolution satellite imagery reveal about economic wealth? *Annual Bank Conference on Development Economics 2016: Data and Development Economics World Bank, September 2016*. Available at documents.worldbank.org/curated/en/610771513691888412/pdf/WPS8284.pdf. Accessed December 20, 2018.
- Jean N, et al. (2016) Combining satellite imagery and machine learning to predict poverty. *Science* 353:790–794.
- Watmough GR, Atkinson PM, Saikia A, Hutton CW (2016) Understanding the evidence base for poverty–environment relationships using remotely sensed satellite data: An example from Assam, India. *World Dev* 78:188–203.
- Okwi PO, et al. (2007) Spatial determinants of poverty in rural Kenya. *Proc Natl Acad Sci USA* 104:16769–16774.
- Tallis H, Kareiva P, Marvier M, Chang A (2008) An ecosystem services framework to support both practical conservation and economic development. *Proc Natl Acad Sci USA* 105:9457–9464.
- Angelsen A, et al. (2014) Environmental income and rural livelihoods: A global-comparative analysis. *World Dev* 64:S12–S28.
- Kristjansson P, et al. (2010) Understanding poverty dynamics in Kenya. *J Int Dev* 22: 978–996.
- Scoones I (2009) Livelihoods perspectives and rural development. *J Peasant Stud* 36: 171–196.
- Zimmerer K, Vanek S (2016) Toward the integrated framework analysis of linkages between agrobiodiversity, livelihood diversification, ecological services and sustainability amid global change. *Land (Basel)* 5:10–38.
- Mamo G, Sjaastad E, Vedeld P (2007) Economic dependence on forest resources: A case from Dendi district, Ethiopia. *For Policy Econ* 9:916–927.
- Stifel D, Minten B (2017) Market access, well-being and nutrition: Evidence from Ethiopia. *World Dev* 90:229–241.
- Cumming GS, Cumming DHM, Redman CL (2006) Scale mismatches in socio-ecological systems: Causes, consequences and solutions. *Ecol Soc* 11:14.
- Mcsweeney K (2002) Who is ‘Forest-Dependent’? Capturing local variation in forest product sale, Eastern Honduras. *Prof Geogr* 54:158–174.
- Jelinski DE, Wu J (1996) The modifiable areal unit problem and implications for landscape ecology. *Landscape Ecol* 11:129–140.
- Seguin A-M, Apparicio P, Riva M (2012) The impact of geographical scale in identifying areas as possible sites for area-based interventions to tackle poverty: The case of Montreal. *Appl Spat Anal Policy* 5:231–251.
- Sanchez P, et al. (2007) The African Millennium villages. *Proc Natl Acad Sci USA* 104: 16775–16780.
- Mutuo P, et al. (2007) Baseline Report: Millennium Research Village Sauri, Kenya, The Earth Institute at Columbia University, p92. Available at mp.convio.net/site/DocServer/Sauri\_Baseline\_Report\_final\_3-7-07.pdf. Accessed date May 9, 2018.
- Kilic T, Serajuddin U, Uematsu H, Yoshida N (2017) Costing household surveys for monitoring progress toward ending extreme poverty and boosting shared prosperity, Policy Research Working Paper, WPS 7951, The World Bank Group, Washington DC.
- Sen A (1999) *Development as Freedom* (Oxford Univ Press, Oxford), 366 p.
- Leach M, Mearns R, Scoones I (1999) Environmental entitlements: Dynamics and institutions in community-based natural resource management. *World Dev* 27:225–247.
- Burke M, Lobell DB (2017) Satellite-based assessment of yield variation and its determinants in smallholder African systems. *Proc Natl Acad Sci USA* 114:2189–2194.
- Wunder S, Börner J, Shively G, Wyman M (2014) Safety nets, gap filling and forests: A global-comparative perspective. *World Dev* 64(Suppl 1):S29–S42.
- Norman P, Pickering C (2017) Using volunteered geographic information to assess park visitation: Comparing three online platforms. *Appl Geogr* 89:163–172.
- Steele JE, et al. (2017) Mapping poverty using mobile phone and satellite data. *J R Soc Interface* 14:20160690.
- Michelson H, Muniz M, DeRose K (2013) Measuring socio-economic status in the Millennium villages: The role of asset index choice. *J Dev Stud* 49:917–935.
- Watmough GR, Sullivan C, Palm CA (2017) An operational framework for object-based land use classification of heterogeneous rural landscapes. *Int J Appl Earth Obs Geoinf* 54:134–144.
- Watmough GR, Atkinson PM, Hutton CW (2013a) Predicting socioeconomic conditions from satellite sensor data in rural developing countries: A case study using female literacy in Assam, India. *Appl Geogr* 44:192–200.
- Watmough GR, Atkinson PM, Hutton CW (2013b) Exploring the links between census and environment using remotely sensed satellite sensor imagery. *J Land Use Sci* 8: 284–30.
- Jönsson P, Eklundh L (2004) TIMESAT—A program for analysing time-series of satellite sensor data. *Comput Geosci* 39:833–845.
- Ripley B (2014) tree: Classification and Regression Trees, R Package, version 1.0-39. Available at <https://cran.r-project.org/web/packages/tree/index.html>. Accessed December 20, 2018.
- James G, Witten D, Hastie T, Tibshirani R (2013) *An Introduction to Statistical Learning: With Applications in R* (Springer, New York), 426 p.
- Kuhn M, Johnson K (2008) Building predictive models in R using the caret package. *J Stat Softw* 28:1–26.