



Continuous and consistent land use/cover change estimates using socio-ecological data

Michael Marshall^{a†}, Michael Norton-Griffiths^a, Harvey Herr^a, Richard Lamprey^b, Justin Sheffield^c, Tor Vagen^a, and Joseph Okotto-Okotto^d

^a Climate Research Unit, World Agroforestry Centre, United Nations Ave, Gigiri, P.O. Box 30677-00100, Nairobi, Kenya, Email: m.marshall@cgiar.org

^b Fauna & Flora International, The David Attenborough Building, Pembroke St, Cambridge, CB2 3QZ, UK

^c Department of Civil and Environmental Engineering, Princeton University, Princeton, NJ, 08544, USA

^d Lake Basin Development Authority, P.O Box 1516-40100, Kisumu, Kenya

[†] Corresponding Author

1

2 **Abstract**

3 A growing body of research shows the importance of land use/cover change (LULCC) on
4 modifying the earth system. Land surface models are used to stimulate land-atmosphere
5 dynamics at the macro- (regional to global) scale, but bias and uncertainty remain that need to be
6 addressed, before the importance of LULCC is fully realized. In this study, we propose a
7 method of improving LULCC estimates for land surface modelling exercises. The method yields
8 continuous (annual) long-term (30-year) estimates of LULCC driven by socio-ecological
9 geospatial predictors available seamlessly across sub-Saharan Africa that can be used for both
10 retrospective and prospective analyses. The method was developed with 2,252 5x5 km² sample
11 frames of the proportion of several land cover types in Kenya over multiple years. Forty-three
12 socio-ecological predictors were evaluated for model development. Machine learning was used
13 for data reduction and simple (functional) relationships defined by generalized additive models
14 were constructed on a subset of the highest ranked predictors ($p \leq 10$) to estimate LULCC. The
15 predictors explained 62% and 65% of the variance in the proportion of agriculture and natural
16 vegetation, respectively, but were less successful at estimating more descriptive land cover types.



1 In each case, population density on an annual basis was the highest ranked predictor. The
2 approach was compared to a commonly used remote sensing classification procedure, given the
3 wide use of such techniques for macro-scale LULCC detection, and out-performed it for each
4 land cover type. The approach was used to demonstrate significant trends in expanding
5 (declining) agricultural (natural vegetation) land cover in Kenya from 1983-2012, with the
6 largest increases (declines) occurring in densely populated high agricultural production zones.

7 **Key words: earth system; land use/cover change; population-environment; remote sensing;**
8 **vegetation function**



1 **1. Introduction**

2 Land use/cover change (LULCC) is an important concern for global environmental
3 sustainability, because it can adversely affect surface albedo and heating (Davin and de Noblet-
4 Ducoudré, 2010); evapotranspiration and other components of the hydrologic cycle (Sterling et
5 al., 2013); local to regional climate with the direct coupling, local coupling, or indirect recycling
6 of surface moisture (Makarieva et al., 2013); global climate via carbon and other greenhouse gas
7 emissions (Anderson-Teixeira and DeLucia, 2011); and ecosystem services worsened by these
8 impacts (Turner et al., 2013). Land surface models, which can be coupled to a regional or global
9 climate model, are used to simulate land-atmosphere interactions retrospectively or prospectively
10 (Pitman, 2003) to identify intervention “hotspots” or develop realistic land management
11 scenarios at the macro- (regional to global) scale (Turner et al., 2007). In general, spatially-
12 explicit LULCC is not an input to land surface models, but is instead represented by structural
13 (e.g. leaf area index) or physiological (e.g. stomatal resistance) changes in vegetation. LULCC
14 is then mapped in parallel to characterize these changes. Even though the importance of LULCC
15 in modifying the earth system has been established and methodologies exist to quantify its
16 impact, studies remain few, due in part to the inadequacy of LULCC estimates (Pielke et al.,
17 2011). In order to further land-atmosphere interaction research, LULCC models must be
18 developed that provide consistent estimates over long (pre-1981) time frames, regular (annual)
19 intervals, and large spatial domains at moderate (5 km) spatial resolution; are projectable 50-100
20 years into the future; and use the same classification approach (Meiyappan et al., 2014;
21 Rounsevell et al., 2014; Verburg et al., 2011).

22 Heistermann et al. (2006) reviews the two primary categories of macro-scale LULCC
23 models (geographic and economic) while Schaldach and Priess (2008) and Rounsevell et al.



1 (2014) include reviews of integrated approaches that combine the strengths of both categories.
2 The Conversion of Land Use and its Effects (CLUE) model (Veldkamp and Fresco, 1996;
3 Verburg et al., 2002) is an example of a geographic technique. It identifies important socio-
4 (population, economy, society, politics and planning, culture, and technology) ecological
5 (climate, vegetation, soil, topography, and hydrology) predictors from observed LULCC data,
6 which are related to each other using multivariate regression or other statistical technique, and
7 then cellular automata are used to simulate competition between the predicted land use/cover
8 types and neighboring grid cells based on these relationships over baseline or projected periods.
9 Decision rules are typically used iteratively to guarantee realistic land use/cover transitions
10 occur. LandSHIFT (Alcamo et al., 2011) is an example of an economic approach, because
11 supply (land use/cover) is distributed on a grid cell basis by demand. Supply is determined from
12 national estimates of crop yield and the net primary productivity of grasslands. Multi-criteria
13 analysis, which involves applying cost functions and land use constraints based on socio-
14 ecological inputs, is used to define demand hierarchically and disaggregate supply over baseline
15 or projected periods. Integrated approaches (e.g. CLUMondo: van Asselen and Verburg, 2013)
16 are becoming more common, because they more adequately account for LULCC processes and
17 the interaction of demand and trade with supply than economic or geographic models,
18 respectively. Like most geographic and economic models, however, integrated models have a
19 sound theoretical basis, but are not widely used for macro-scale applications, because of data
20 inconsistencies and incongruities and model complexity that can propagate error, as well as, the
21 time and other resources needed to operate them. Earth observation (remote sensing) models are
22 an important sub-category of the geographic approach, because they overcome many of these
23 challenges, making their operational use at the macro-scale more feasible.



1 Hansen and Loveland (2012) and Ban et al., (2015) present recent reviews of macro-scale
2 remote sensing-based LULCC modeling. Remote sensing approaches use multivariate statistical
3 techniques to classify land cover types based on the spectral or textural characteristics of gridded
4 satellite data (DeFries et al., 1995). These approaches are simpler than integrated approaches,
5 because they tend to capture change at a single resolution directly with no interaction between
6 adjacent pixels. Remote sensing approaches, therefore, have the potential to reach higher model
7 parsimony than integrated approaches, and require less time for processing and considerably
8 fewer data types. Early remote sensing approaches involved daily coarse spatial resolution
9 (8km) Advanced Very High Resolution Radiometer (AVHRR) data available from 1981. Large
10 disagreement and uncertainties in the models, due to mixed pixel effects from small land
11 use/cover patch size, as well diverse classification systems and methods, have limited their use at
12 the macro-scale (Lepers et al., 2005). Improved computational storage and processing and
13 consensus on classification has facilitated the creation of consistent global LULCC maps at
14 Landsat (30m) resolution (Giri et al., 2013). GlobeLand30 (Chen et al., 2015), for example, uses
15 a pixel-object-knowledge-based approach to classify Landsat images from spectrally-derived
16 vegetation indices globally in 2000 and 2010. The use of Landsat data alone poses serious
17 challenges to modeling LULCC on an annual basis: persistent cloud cover, particularly in the
18 tropics during the primary growing season, and a 16-day revisit cycle, makes retrieval of cloud-
19 free pixels difficult; the Landsat platforms have been retired (Landsat 5), have failed (Landsat 6),
20 suffer from technical problems (Landsat 7), or are only recently active (Landsat 8). To improve
21 the temporal resolution and continuity of classification, other products, such as the Global Forest
22 Change product (Hansen et al., 2010), fuse Moderate-resolution Imaging Spectroradiometer
23 (MODIS) data available every 1-2 days at 250-500m spatial resolution with Landsat data. But



1 these products are only available over the MODIS era (2000-present), making long-term
2 classification infeasible. In short, the major drawback of remote sensing approaches is that the
3 temporal range and continuity necessary for long-term annual global change detection are often
4 sacrificed for high ($\leq 500\text{m}$) spatial resolution. Finally, remote sensing data is not projected 50-
5 100 years into the future or available pre-1981 like other socio-ecological data, such as
6 population density, precipitation, or temperature, limiting their use for retrospective or
7 prospective analysis.

8 The purpose of this study was to propose a simple (functional) way to map LULCC at the
9 macro-scale at 5 km resolution on an annual basis using socio-ecological predictors that are
10 available pre-1981 and projected 50-100 years into the future in order to facilitate land-
11 atmosphere modeling and research. The method was developed using sample frames consisting
12 of continuous land cover developed from multi-year aerial and ground surveys in Kenya over a
13 30-year period and socio-ecological predictors that are available seamlessly across sub-Saharan
14 Africa (SSA). The approach is compared with remote sensing predictors that have been used to
15 classify land cover types based on their unique phenology. Kenya is an ideal location to develop
16 such a method, because like many countries in SSA, data is scarce compared to the Global
17 North, and the impact of land modification on people and the environment is high (Lambin et al.,
18 2003). In addition: 1) population density is highest in the most agriculturally productive areas
19 due to unequitable land distribution and poor infrastructure (Jayne and Muyanga, 2012), making
20 ecological determinants that are primarily used to map LULCC potentially less relevant (Pricope
21 et al., 2013); 2) agriculture is the primary source of livelihood and crops are mostly rainfed
22 (Ngetich et al., 2014); and 3) inter-annual rainfall variability is high and frequently causes
23 devastating droughts and floods (Held and Soden, 2006).



1 **2. Data and Methods**

2 *2.1 Study area*

3 Aerial surveys were conducted in 1983, 1985, 2012, and 2013, to assess changes in land
4 cover over parts of the Lake Victoria basin and central region of Kenya (Machakos and Makueni
5 areas). The surveys yielded 2,252 5x5 km² grid frames covering 28,150 km² or approximately
6 47% of Kenya's arable lands (**Figure 1**). Olofsson et al. (2012) has suggested that 5x5 km²
7 frames are appropriate for macro-scale LULCC analyses. The lakeshore and lowlands of Lake
8 Victoria basin are primarily tropical with one long rain season that extends from February to
9 September (UNEP, 2008). The neighboring highlands follow a bimodal pattern and annual totals
10 are higher than near the lakeshore, due to warm moist westerlies during the West African
11 monsoon and orographic uplift. Central Kenya is drier and has two distinct rain seasons: long
12 rains (March-June) and short monsoon rains (October- December). The Machakos area, which
13 includes Muranga', Kiambu, and the northern part of Machakos, is humid subtropical and
14 therefore wetter than Makueni to the south-east, which is semi-arid.

15 For each frame, the probability (proportion) of various land cover types was classified at
16 two levels of specificity: level one (agriculture, natural vegetation, urban, and miscellaneous)
17 and level two (crops, fallow, shrubs, savanna, wetlands, forest, and agroforestry). Continuous
18 data were used, because at 5 km resolution, spatial heterogeneity makes discrete classification
19 impractical. Agriculture included agroforestry, defined here as trees on a farm; crops (banana,
20 coffee, maize, sugar cane, tea, wheat, and others); and pasture/fallow. Natural vegetation
21 included savanna, shrubs (open and closed), wetlands (perennial and permanent), and forest
22 (evergreen and deciduous). Urban included built up structures, such as roads, homes, and towns.
23 Miscellaneous included fish ponds and other water bodies, exposed rock, and charcoal pits. The



1 frames were developed using an aerial point sampling approach (Norton-Griffiths, 1988): several
2 thousand geotagged aerial photos were taken over parallel transects spaced 1 km apart at
3 approximately 488 m (height-above-ground) in 1983/1985 and then again in 2012/2013,
4 resulting in approximately 7 aerial natural color analogue photos per frame with a ground-
5 sampling-distance of < 1 cm in 1983/1985 and 5 aerial natural color digital photos per frame
6 with a ground-sampling-distance of 6.5 cm in 2012/2013. The retrieval dates are shown in
7 **Table 1**. A team of six technicians interpreted the photos on a rolling basis to minimize potential
8 bias and errors that can occur from manual classification by different interpreters and for
9 different years. The proportion of each land cover type (0-100%) was determined by manually
10 classifying a grid of 320 randomly distributed points superimposed over each photo. For each
11 year, all land cover types were represented and classified, but not all sample frames were
12 interpreted and classified (**Figure 1**). The interpretations were validated via site visits and
13 meetings with community stakeholders. The estimates were then averaged over the photos
14 across interpreters to get the proportions for each frame. Further details on the 1983/1985 and
15 2012/2013 campaigns can be found in EcoSystems Ltd (1983), EcoSystems Ltd (1987), and
16 Lamprey (2013).

17 *2.2 Macro-scale data handling and processing*

18 Forty-three non-remote sensing (climatic, hydrologic, socioeconomic, and topographic)
19 and sixteen remote sensing (phenological) predictors of land cover change were compared and
20 subset for model-building with the sample frames. The predictors were selected, because they
21 are gridded seamlessly across SSA and could therefore facilitate continuous and consistent land
22 cover classification across the continent. Either slowly-changing (long-term average/one-time
23 value) or annually-changing predictors were considered. The slowly-changing predictors and



1 their sources are shown in **Table 2**. Using these predictors alone could streamline the modeling
2 process. However, in reality, phenology, climate, and population change frequently, so these
3 predictors were derived on an annual basis as well. The handling and processing of annually-
4 changing predictors are discussed in Sections 2.2.1-2.2.3. For the remainder of the paper,
5 annually-changing variables include a “.d” extension. All of the data was projected to Africa
6 Equidistant Conic (m) to facilitate distance calculations. The predictors were resampled to the
7 finest resolution data (90 m) and aggregated to 5 km resolution for model-building.

8 2.2.1 Climate

9 BIOCLIM variables were chosen for the analysis, because they 1) are commonly used for
10 similar studies that require biologically meaningful climate information and 2) have been
11 projected mid-21st century at high spatial resolution for SSA (see Platts et al., 2014). Two
12 additional climate parameters were included in the analysis, because they are part of the Platts et
13 al. (2014) dataset: atmospheric demand for moisture (Potential Evapotranspiration- PET) and
14 the Moisture Index. The BIOCLIM variables were computed on an annual basis from 1983-
15 2012 using monthly temperature, shortwave incoming radiation, and precipitation. The variables
16 were computed using the “biovars” function in the “dismo” package in R (Hijmans et al., 2015).
17 As with the Platts et al. (2014) dataset, PET was estimated using the Hargreaves and Samani
18 (1985) approach.

19 The temperature/radiation and precipitation predictors were taken from the Princeton
20 University high resolution meteorological forcing (PHF) (Chaney et al., 2014) and the Climate
21 Hazards Group InfraRed Precipitation with Stations (CHIRPS) (Funk et al., 2014) datasets,
22 respectively. PHF originally spanned 1979-2008, but was extended to 2012 for this study. It is a
23 downscaled version of the Princeton University global meteorological forcing (PGF) dataset



1 (Sheffield et al., 2006) for SSA. It assimilates new observation data, specifically station data
2 from the U.S. National Climatic Data Center (NCDC) Integrated Surface Database (ISD) and has
3 undergone more rigorous correction than the global dataset. PHF is a blend of the most up-to-
4 date observation-based, remote sensing, and reanalysis data sources: the National Centers for
5 Environmental Prediction–National Center for Atmospheric Research (NCEP-NCAR) reanalysis,
6 Global Precipitation Climatology Project, Tropical Rainfall Measuring Mission (TRMM), the
7 Climatic Research Unit (CRU), and the Surface Radiation Budget. Downscaling is performed
8 using bilinear interpolation weighted by elevation. The dataset includes precipitation,
9 minimum/maximum temperature, pressure, shortwave and longwave radiation, specific
10 humidity, and wind speed at a daily time step and 0.1° (~10 km at the equator) resolution.
11 CHIRPS is available at pentad (5-day) intervals and 0.05° (~5km at the equator) spatial
12 resolution from 1981-2012. Like PHF, CHIRPS is a blend of several observation-based, remote
13 sensing, and reanalysis sources: geostationary thermal infrared satellite observations from the
14 Climate Prediction Center and National Climatic Data Center; TRMM; and NOAA-NCAR.
15 CHIRPS was selected as the precipitation source over PHF, because it incorporates the largest
16 collection of ground-based precipitation data in East Africa and bias-correction is performed
17 using the Climate Hazards Precipitation Climatology (Funk et al., 2015).

18 2.2.2 Population density

19 Population density was derived from the UNEP/GRID-Sioux Falls African Population
20 Distribution Database (APDD) on an annual basis from 1983-2012. APDD consists of
21 population density at a spatial resolution of 2.5 arc-minutes $^\circ$ (~5km at the equator) for base
22 years 1960, 1970, 1980, 1990, and 2000. The grids are derived from population statistics at
23 various administrative levels (defined by vector polygons) and temporal scales, depending on the



1 availability of national population statistics. The approach taken to convert population polygons
2 to gridded population is detailed in Deichmann (1996). Each grid cell represents “population
3 potential”, based on its proximity to the transportation network (roads, railroads, and navigable
4 rivers, and major towns/cities). Population at a given administrative level is then disaggregated
5 according to the population potential. Grid cells that are closer to the network have higher
6 coefficients and therefore receive a larger proportion of the population than grid cells further
7 away. The base years are then extrapolated with an exponential growth/decay function (Davis,
8 1995). For consistency, the same function was used to distribute population between base years
9 on an annual basis for each grid cell:

$$P_{i,j,t} = P_{i,j,T} e^{\Delta t k_{i,j}} \quad (2)$$

$$k_{i,j} = \ln(P_{T+10n}/P_{T+10(n-1)})/10 \quad (3)$$

10 $P_{i,j,t}$ is the interpolated population/population density for a given year (t) and at grid cell i, j , $P_{i,j,T}$
11 is the population/population density for a given base year (period = 10 years), Δt is the change in
12 time from the base year to the year being interpolated, and $k_{i,j}$ (**Equation 2**) is the growth/decay
13 coefficient. The growth/decay coefficient is defined by $P_{T+10(n-1)}$ (initial base year for iteration n)
14 and P_{T+10n} (last base year for iteration n). The denominator was set to ten, because $k_{i,j}$ accounted
15 for decadal trends. After 2000, population statistics were extrapolated to 2012 using the 1990-
16 2000 growth/decay coefficients.

17 2.2.3 Remote Sensing Predictors

18 The National Aeronautics and Space Administration’s Global Inventory Modeling and
19 Mapping Studies (GIMMS) Normalized Difference Vegetation Index (NDVI) Version 3
20 (NDVI3g) (Pinzon and Tucker, 2014) was used to estimate the remote sensing predictors. NDVI
21 is a ratio-based vegetation index derived from Earth observation (AVHRR) surface reflectance in



1 the visible red and near infrared (NIR). NDVI approaching one (zero) is indicative of dense
2 vegetation (bare soil). NDVI3g is available at 0.08° (~8km resolution at the equator) spatial
3 resolution and at a 15-day timestep from 1983-2013. NDVI3g has been compared to other long-
4 term global vegetation records and is considered the most appropriate for trend analysis (Tian et
5 al., 2015).

6 The predictors were derived from NDVI using harmonic regression (Eastman et al.,
7 2009) on an annual basis from 1983-2012. Linear harmonic regression estimates the amplitude
8 (maximum) and phase (timing) of a fitted time series, but unless higher order harmonics are
9 introduced, linear harmonic regression is susceptible to outliers and multimodal regimes
10 commonly found in the tropics. To overcome these obstacles, non-linear harmonic regression
11 (Carrão et al., 2010) was used to estimate five phenological predictors:

$$\mathbf{NDVI}_{i,j,T} = \mathbf{M}_{i,j} + \mathbf{A}_{i,j} \cos(\omega_0 t + \phi + \alpha \cos(\omega_0 t + \varphi)) \quad (1)$$

12 Where $\mathbf{NDVI}_{i,j,T}$ is NDVI3g at grid cell i, j and over period T , which in this case was 24,
13 because non-linear harmonic regression was computed on an annual basis from the 15-day data;
14 \mathbf{M} is the intercept (annual mean NDVI); \mathbf{A} is the amplitude; ϕ is the annual phase; and α and φ
15 are non-linear terms defining the strength of non-linearity (asymmetry) and non-linear phase
16 (deceleration/acceleration of asymmetry), respectively. The frequency (ω_0) equals $2\pi/T$. The
17 approach can be reduced to a linear harmonic oscillator by setting $\alpha \cos(\omega_0 t + \varphi)$ to zero. The
18 non-linear predictors were derived at each grid cell using the “nlsLM” function in the
19 “minipack.lm” package in R (Elzhov et al., 2015). nlsLM uses the Levenberg-Marquardt
20 optimization method (Moré, 1978) to find the non-linear least-squares fit. The function was
21 constrained by the seed and boundary conditions described in Carrão et al., (2010). One
22 thousand iterations at each grid cell were performed to avoid fitting local optima. Linear terms



1 (A and ϕ) were computed for the analysis as well, using the “lm” function in the “stats” package
2 in R (<https://cran.r-project.org/>), because they are more efficient and are easier to interpret.

3 *2.3 Land-cover model development using remote sensing and non-remote sensing predictors*

4 Land cover models were developed at both levels of specificity and involved three steps:
5 1) data reduction and model feasibility; 2) functionalizing the relationships between selected
6 predictors and each land cover class; and 3) evaluation. Seventy percent of the samples
7 (N=1,576) were used for model calibration and 30% of the samples (N=676) were used for
8 model validation.

9 Machine learning was used to omit redundant predictors and determine the feasibility of
10 using the remaining predictors to predict each land cover type, given the large number of
11 predictors and possible inter-correlations. Machine learning techniques lead to stable results
12 when the number of predictors is large and are less affected by non-linearity and
13 multicollinearity than other automated fitting routines (Binder and Tutz, 2008). The Breiman's
14 random forest algorithm (Breiman, 2001) available in the “randomForest” package in R was
15 selected in particular, because it is less susceptible to over-fitting and yields higher prediction
16 accuracy than other machine learning algorithms (Fernández-Delgado et al., 2014). The Random
17 Forest (RF) algorithm yields an ensemble model, bagged from multiple and independent decision
18 trees consisting of various combinations of predictors and sample subsets. The performance of
19 the ensemble is measured with a pseudo coefficient of determination (pseudo- R^2), which is one
20 minus the ratio of the cross-validated mean squared error (MSE) of the prediction to the variance
21 of the observed data. The importance of each predictor in the ensemble is also quantified and is
22 defined by the percent increase in cross-validated MSE when a predictor is removed from the
23 ensemble. Once the predictors were ranked, the “rfcv” function was used to determine the



1 number of predictors to use to develop functional relationships for each land cover class. Rfcv
2 computes the cross-validated MSE versus the number of predictors included in the ensemble in
3 descending order of importance.

4 The drawback of RF is that it results in complex relationships that are difficult to
5 interpret. Generalized Additive Models (GAMs) (Hastie and Tibshirani, 1990) were used to
6 build functional relationships on the subsets of important predictors identified with RF, because a
7 number of studies have successfully estimated the proportion of crop area with socio-ecological
8 predictors and GAMs (Grace et al., 2014; Husak et al., 2008; Marshall et al., 2011); like RF,
9 GAMs are not severely impacted by non-linear data; and unlike RF, GAMs are relatively simple
10 and easy to interpret. Since the response variable (proportion of land cover type) was continuous
11 and bounded from 0-100%, the data was fitted using a quasi-binomial distribution
12 (link=logistic). The logistic GAM predicts the log-likelihood of an event (probability of
13 success/probability of failure) using, in our case, a series of cubic spline functions:

$$\log\left(\frac{p_j}{1-p_j}\right) = \beta_0 + \sum f_{i,j}(x_{i,j}) \quad (4)$$

14 Where \mathbf{p} is the probability of a land use/cover type for sample frame \mathbf{j} , β_0 is the intercept, and
15 $f_{i,j}(x_{i,j})$ is the cubic spline function for predictor x_i at sample frame \mathbf{j} . The GAMs were developed
16 with the “gam” function in the “mgcv” package in R. Model calibration was evaluated with
17 explained part and overall deviance. Deviance is the log-likelihood alternative to variance. Part
18 deviance is the deviance explained when the target predictor is removed from a GAM minus the
19 overall deviance. Another pseudo- R^2 statistic (1-model deviance/null deviance) was also
20 computed to compare calibration statistics with validation statistics (R^2 and MSE).

21 In order to demonstrate how the models can be used for macro-scale application, the final
22 GAMs developed were used to reconstruct the annual change in agriculture and natural



1 vegetation and to perform a trend analysis from 1983-2012 at each sample frame. Trends were
2 estimated using the Theil-Sen technique, which computes the median of all possible pairwise
3 slopes in a time series. The approach has been used, for example, to measure long-term trends in
4 NDVI (de Beurs and Henebry, 2005), because it is not significantly impacted by outliers or non-
5 linearity. The significance of each trend was assessed using the Mann-Kendall statistic. Trends
6 were masked at the 99.9% confidence band.

7 **3. Results**

8 *3.1 Land cover sample frame summary*

9 The distribution of land cover over the sample frames is illustrated with a boxplot in
10 **Figure 2**. Agriculture and natural vegetation land cover (level one) were normally distributed,
11 with agriculture having a higher median (54.04%) and lower spread (29.32% and 76.33% at the
12 first and third quartiles) than natural vegetation (median=39.72%, first quartile=16.21%, and
13 third quartile=65.67%). The proportion of urban and miscellaneous land cover was considerably
14 lower (median=4.00% and 0%, respectively) and non-linear, each having several high proportion
15 outliers. The disaggregated land cover (level two) distributions, with the exception of crops,
16 were non-linear, each having long right tails. Crops represented the largest proportion of land
17 cover (median=37.52%) and had the largest spread (19.43% and 58.46% at the first and third
18 quartiles), followed by savanna (median=15.79%, first quartile=3.80%, and third
19 quartile=31.91%). Wetlands represented the smallest proportion of land cover (median=0%),
20 with sample frames not exceeding 75%, while forest represented the second smallest proportion
21 of land cover (median=2.22%), but had the longest right tail with proportions reaching 100%.

22 *3.2 Data reduction*



1 The top remote sensing and non-remote sensing predictors considered are ranked in
2 descending order of importance for agriculture and natural vegetation using bar graphs in **Figure**
3 **3**. The RF ensemble models using non-remote sensing predictors performed moderately well for
4 agriculture (pseudo- $R^2=0.69$) and natural vegetation (pseudo- $R^2=0.69$), but poorly for the more
5 non-linear land cover distributions (urban pseudo- $R^2=0.37$ and miscellaneous pseudo- $R^2=0.50$).
6 The RF ensemble models using remote sensing predictors all performed poorly, but incorporated
7 a smaller number of predictors than the non-remote sensing ensembles: agriculture (pseudo-
8 $R^2=0.49$), natural vegetation (pseudo- $R^2=0.50$), urban (pseudo- $R^2=0.22$), and miscellaneous
9 (pseudo- $R^2=0.33$). For the non-remote sensing ensembles, annually-changing predictors were
10 more important than slowly-changing predictors, and population density and climate predictors
11 consistently outranked topographic or hydrologic predictors. Popd.d, popd, bio7.d, bio14.d, and
12 bio3.d were consistently ranked the most important predictors of agriculture and natural
13 vegetation proportions. Omitting popd.d, the most important predictor for agriculture, for
14 example, led to a more than 65% increase in ensemble MSE. Given that popd.d and popd were
15 both important, model results were compared with popd.d and popd individually and combined
16 as anomalies (popd.d/popd). Ensemble performance was better when the two predictors were
17 considered separately. The most important remote sensing predictors were less influential than
18 popd.d. strn, ampn.d, and ampl.d were more equally important for agriculture and natural
19 vegetation, followed by phsl and phsn.

20 The importance of predictors of level two (crops, savanna, and forest) proportions are
21 ranked in **Figure 4**. The ranking was more variable for level two classifications, but popd.d
22 remained the most important predictor in each case. The level two RF ensemble models
23 predicted less variability than the level one RF ensemble models and the non-remote sensing



1 predictors outperformed the remote sensing predictors. The non-remote sensing models
2 performed moderately well for crops (pseudo- $R^2=0.63$), savanna (pseudo- $R^2=0.62$), and forest
3 (pseudo- $R^2=0.61$), but poorly for fallow (pseudo- $R^2=0.42$), shrubs (pseudo- $R^2=0.54$), wetlands
4 (pseudo- $R^2=0.10$), and agroforestry (pseudo- $R^2=0.55$). Precipitation-based climatic predictors
5 (bio12.d, bio13.d, bio14.d, and bio16.d) were more important in the savanna ensemble than
6 temperature-based climatic variables driving the crop ensemble. For the forest simulation,
7 topographic predictors (slp and topind) were more important than most of the climatic predictors.
8 The remote sensing ensembles performed poorly for all of the level two land cover classes: crops
9 (pseudo- $R^2=0.46$), fallow (pseudo- $R^2=0.33$), shrubs (pseudo- $R^2=0.44$), savanna (pseudo-
10 $R^2=0.44$), wetlands (pseudo- $R^2<1\%$), forest (pseudo- $R^2=0.46$), and agroforestry (pseudo-
11 $R^2=0.41$). For crops, strn and ampl.d remained the most important predictors. Maximum annual
12 NDVI, as captured by ampl.d and ampn.d, were much more important for predicting the
13 proportion of savanna. Unlike other ensembles, which were driven by annually-changing
14 predictors, the most important remote sensing predictors for forest cover were long-term
15 averages.

16 *3.3 Building functional relationships*

17 The GAMs were developed for moderately performing land cover classes and used
18 considerably fewer predictors than the RF ensembles, because most of the predictors in the
19 ensembles explained very little, if any variability. This is illustrated in **Figure 5**, which shows
20 MSE versus the number of predictors used in the non-remote sensing and remote sensing
21 ensembles for forest. For the non-remote sensing ensemble, MSE increased from 119.76 to
22 120.49 after the 10th predictor and leveled off after the 13th predictor was introduced. For the
23 remote sensing ensemble, MSE increased from 120.49 to 163.34 and levelled off after the 7th



1 predictor was introduced. For this reason, the GAMs were built with 10-13 of the highest ranked
2 non-remote sensing predictors and additional predictors, namely popd, were removed after
3 redundancies were identified in the GAM component functional plots and with significance tests
4 (not shown). GAMs were not constructed using the remote sensing predictors, because of the
5 poor results of the ensembles and the inability of additional predictors to substantially improve
6 the accuracy of the GAMs. Similarly, non-remote sensing GAMs were not developed for urban,
7 miscellaneous, fallow, shrubs, or wetlands.

8 **Figures 6 and 7** show the functional relationships of the predictors used for estimating
9 the proportion of agriculture and natural vegetation. Each model explained 61.5% (pseudo-
10 $R^2=0.66$) and 61.4% (pseudo- $R^2=0.66$) of model deviance with nine and seven predictors,
11 respectively. The confidence intervals tended to be wider at proportion extremes, because fewer
12 data points were available to train the models. The relative importance of each predictor, as
13 defined by part deviance and other calibration statistics are shown in **Table 3** for the land cover
14 types that were considered feasible for model-building. Popd.d remained the most important
15 predictor and uniquely explained 7.0-26.2% of model deviance. The log-likelihood of
16 agriculture (natural vegetation) increased (decreased) rapidly as population density increased
17 from 0 to 550 people•km⁻², more gradually between 550 and 1200 people•km⁻², and reversed
18 beyond 1200 people•km⁻². The predictive power of the topographic and climatic variables
19 dropped off sharply after popd.d. For agriculture, bio14.d and topind were the second and third
20 most important predictors, but explained only 1.9% and 1.6% unique deviance. As seen in the
21 partial functional plots, the proportion of agriculture was highest in high production zones
22 (medium population density) on ridges and crests where topind was low and for very wet tropical
23 areas where bio14.d was high and semi-arid areas where bio14.d was low. For natural



1 vegetation, temperature predictors, bio4.d and bio7.d, explained the second and third highest
2 unique deviance after popd.d (2.0% and 1.3%). As seen in the functional plots, low populated
3 areas with more temperature seasonality, or inter-annual variation, and less isothermality (higher
4 non-tropical latitudes) tended to have higher proportions of natural vegetation (savanna and
5 shrubs). For the level two classifications, calibration was more difficult and yielded poorer
6 relationships. Popd.d was the most important predictor and explained 7.0 – 16.4% unique
7 deviance. The predictive power of the topographic and climatic variables was more equally
8 distributed than for the level one classification.

9 In all cases, the R^2 for the validation subset was lower than the pseudo- R^2 from the
10 calibration subset: agriculture ($\Delta R^2 = -0.04$), natural vegetation ($\Delta R^2 = -0.01$), crops ($\Delta R^2 = -0.03$),
11 savanna ($\Delta R^2 = -0.01$), and forest ($\Delta R^2 = -0.06$) (**Figure 8**). With the exception of the crops GAM,
12 level two GAMs tended to under-predict high proportions of land cover (savanna and forest) and
13 contained numerous outliers.

14 *3.4 Trend analysis*

15 The GAMs for agriculture and natural vegetation were used to simulate trends in the
16 annual proportions for the sample frames from 1983-2012 as part of the evaluation to
17 demonstrate how the approach could be used for a retrospective analysis. The proportion of
18 agriculture for 1983 and 2012 are shown in **Figure 9a and 9b**, while trends over the 30 year
19 period are shown in **Figure 9c**. The high potential agricultural zone (wet highlands) in Western
20 Kenya experienced the largest increase in simulated agricultural cover ($> 1\%$ per year or 30%
21 over the 30-year period). A time series of the strongest trend (1.68% per year) is shown in
22 **Figure 9d**. Simulated population density was at 145 people \cdot km $^{-2}$ in 1983 for this sample frame,
23 which steadily increased to 478 people \cdot km $^{-2}$ in 2012. Closer to the lake, which consists of drier



1 marginal mixed farming, trends were insignificant at the 99.9% confidence band or relatively
2 weak ($< 1\%$ per year). Similar patterns were seen for the marginal mixed farming and high
3 potential agricultural zones of central Kenya as well. The only decreasing trend in agricultural
4 cover was seen in Kitale town (-1.40% per year). The time series is also shown in **Figure 9d**.
5 Population growth in Kitale was $1,110 \text{ people}\cdot\text{km}^{-2}$ in 1983, which is near the threshold of
6 declining agriculture cover versus population density at $1,200 \text{ people}\cdot\text{km}^{-2}$. By 2009, when the
7 largest decrease in agriculture cover occurred, from 51.0 to 29.5%, population density had
8 steadily increased and surpassed another apparent threshold above $3,000 \text{ people}\cdot\text{km}^{-2}$. The
9 direction and relative magnitude of trends in natural vegetation (not shown) generally
10 corresponded inversely to trends in agriculture, but were negatively-weak (maximum= -0.4% per
11 year or -12% over the 30-year period).

12 **5. Discussion**

13 The results make two important contributions that the land surface modeling community
14 should consider to improve LULCC detection: 1) a socioeconomic variable (population density)
15 was the highest ranked predictor of LULCC and had considerably more predictive power than
16 ecological predictors and 2) non-remote sensing predictors in all cases out-performed remote
17 sensing predictors.

18 The global increase in agricultural land cover has been attributed to the demand for food
19 and other agricultural commodities by a growing population (Pongratz et al., 2008). In SSA,
20 smallholder farms, which support the majority of the labor force, are small (half are $< 1.5 \text{ ha}$) and
21 concentrated in densely populated areas, while large portions of arable farmland in sparsely
22 populated areas remain underutilized (Jayne et al., 2003). This underutilization is due primarily
23 to a lack of investment in infrastructure and unequitable tenure systems, which forces farmers to



1 grow more on less land. This relationship is confirmed by rural population survey data in Kenya,
2 which showed that fertilizer input use and net farm income per hectare increase until
3 approximately $550 \text{ persons}\cdot\text{km}^{-2}$ and then sharply decline, because farm sizes shrink, surplus
4 production decreases, and farmers must adopt costlier strategies (e.g. zero-grazing) to maximize
5 revenue (Jayne and Muyanga, 2012). The functional relationship for population density and
6 steady increase in area under cultivation in high production zones demonstrated by the trend
7 analysis in this study, corresponds to this finding, as area under cultivation increased rapidly to
8 approximately $550 \text{ persons}\cdot\text{km}^{-2}$ and then increased more gradually with higher population
9 density until $1200 \text{ persons}\cdot\text{km}^{-2}$. Few sample frames had population densities greater than $1,200$
10 $\text{persons}\cdot\text{km}^{-2}$, as was seen in Kitale town, so it is difficult to know if this functional relationship
11 holds for very high population densities. At least to 2008, Kitale experienced a growth rate of
12 12%, well above the national average (7%), due to persistent drought and out-migration from
13 neighboring high production zones (Majale, 2008), so perhaps this relationship is reasonable.
14 Although the functional relationship for population density corroborates household surveys in
15 Kenya and other agrarian countries in SSA, it should be further scrutinized, because land tenure
16 in SSA is complex (Place, 2009) and the dependency of LULCC predictors on location and
17 spatial scale can be high (Rindfuss et al., 2004).

18 Population density estimates vary widely (Wilson, 2014) and given its fundamental
19 importance to the proposed model framework, future work should aim to integrate a more
20 dynamic product that better accounts for inter-annual variability and realistic representation of
21 current and projected population density. To the authors' knowledge, this was the first attempt
22 to make a population product dynamic (annually-changing). However, the approach is
23 essentially tracking decadal trends that explain a significant portion of inter-annual variability.



1 In reality, population density can show high inter-annual variability due to migration and other
2 factors. Regarding the product itself, changes in population density do not necessarily “grow”
3 from transportation networks and are influenced by important feedbacks. In addition, the
4 extrapolation method used is efficient and can be projected indefinitely, but does not capture
5 complex demographics that other methods do and can lead to “runaway” growth/decline (Baker
6 et al., 2008). Finally, there is no consensus on which population product to use, however, in the
7 future, other products (e.g. Afripop) should be compared against the product used here or
8 combined to make a model ensemble.

9 This paper highlights the importance of gridded socioeconomic data in mapping LULCC,
10 but gridded macro-scale datasets are almost exclusively ecological in nature. The biggest gains
11 in LULCC prediction could be made, therefore, by developing gridded macro-scale
12 socioeconomic data from existing country-level products, such as the Human Development
13 Index. More minor gains could be made by integrating ecological predictors not used in this
14 study, such as soil type and properties. Gridded soils data exists globally from the International
15 Soil Reference and Information Center, but was not considered in this study, because it is a one-
16 time value and does not capture the dynamic nature of soils or its complex relationship with
17 LULCC. A dynamic soils product was recently developed for the MODIS era (see Vågen et al.,
18 2016) and could be a powerful tool for LULCC detection, especially if it is back-casted over the
19 full temporal range of other predictors.

20 Grace et al. (2014) developed GAMs to predict cropped area in Kenya using ecological
21 predictors (rainfall, elevation, NDVI, slope, and the topographic wetness index) and explained
22 much of the deviance in cropped area (41.9-81.4%). Although the models used different
23 predictors for different years and production zones, and the definition of cropped area and the



1 degree of smoothing were not explicit, the study highlights that multicollinearity may be
2 obscuring the importance of ecological predictors. Population density tends to be highly
3 correlated with and could be suppressing the explanatory power of these predictors, though the
4 partial deviance statistics did not reflect this. In addition, the random forest algorithm accounts
5 for multicollinearity, but other techniques could be introduced to further reduce these effects.
6 For example, Principal Components Analysis could be used to develop temperature and
7 precipitation indices that integrate all or some of the BIOCLIM predictors, given the large
8 number analyzed.

9 Phenological patterns extracted from continuous Earth observation based NDVI have
10 been widely used to map LULCC over long time periods, given the lack of higher spatial and
11 spectral resolution data before the MODIS era (Ali et al., 2014; Bie et al., 2012). These studies
12 show that vegetation periodicity is highly variable for a given land cover type and that long-term
13 averages of phenological predictors are more reliable for mapping land use/cover. Indeed, many
14 of the most important remote sensing predictors (particularly for forests) were long-term
15 averages, but they still greatly under-predicted LULCC when compared against non-remote
16 sensing predictors. Perhaps the main difficulty in using long-term Earth observation data for
17 LULCC estimation is the coarseness of the data and the rapid change in vegetation that often
18 occurs over small spatial scales. Population density, which was a much stronger predictor, on
19 the other hand, may well be captured using moderate resolution data, because this predictor
20 changes more gradually over space. An analysis of the non-remote sensing and remote sensing
21 predictors together (not shown) revealed a potential avenue to improve remote sensing LULCC
22 detection using long-term vegetation records. The combination of non-remote sensing and
23 remote sensing predictors over the period analyzed moderately increased the accuracy of



1 estimates for natural vegetation, savanna, and forest land cover types. Meaning for these
2 important land cover types, large gains could potentially be made by integrating long-term
3 vegetation datasets downscaled with Landsat imagery into the model framework.

4 **6. Conclusion**

5 This study developed and evaluated a simple method to provide consistent estimates of
6 LULCC annually over 30 years at 5 km resolution using non-parametric functional relationships
7 with a small subset of socio-ecological predictors ($p \leq 10$). Functional relationships were
8 developed after data mining 43 geospatial datasets that are available seamlessly across SSA,
9 which can be used for retrospective pre-1981 or prospective mid- and late-21st century analyses.
10 The relationships are intuitive and tunable, making their use practical for decision-makers to
11 identify intervention hotspots and develop land management scenarios. Model validation,
12 performed with multi-temporal proportions of major land cover types in Kenya, revealed that a
13 number of activities should be performed to improve the predictive power of the models for
14 practical use. These activities primarily focus on integrating improved existing or newly
15 developed geospatial (particularly socioeconomic) datasets into the proposed model framework.
16 With these improvements, land surface and LULCC modelling could be greatly enhanced and
17 the consequence of the latter on the earth system can be more fully understood.



1 **Acknowledgements**

2 The work summarized in this manuscript was primarily funded through support by the
3 CGIAR research program on Climate Change, Agriculture and Food Security for the project,
4 titled “Multi-disciplinary species distribution modelling for “climate smart” agriculture in East
5 Africa.” Additional support for the early field and aerial surveys was supported by the Kenyan
6 Lake Basin Development Authority and Ministry of Planning and National Development. We
7 would like to extend our special thanks to Sandra Nakibilango, Dorcas Ninsiima, Charles Ngugi,
8 Patrick Ojorot and Samuel Olowo who interpreted the aerial photos taken for this project;
9 Bernadette Apio who coordinated the interpreter team; and Juliet Kyakobyewo who entered the
10 data into our database. Finally, we would like to thank Dr. Eike Luedeling, who initially led and
11 organized the project.



References

- 1 Alcamo, J., Schaldach, R., Koch, J., Kölking, C., Lapola, D., Priess, J., 2011. Evaluation of an
2 integrated land use change model including a scenario analysis of land use change for
3 continental Africa. *Environ. Model. Softw.* 26, 1017–1027.
- 4 Ali, A., de Bie, C.A.J.M., Skidmore, A.K., Scarrott, R.G., Lymberakis, P., 2014. Mapping the
5 heterogeneity of natural and semi-natural landscapes. *Int. J. Appl. Earth Obs.*
6 *Geoinformation* 26, 176–183.
- 7 Anderson-Teixeira, K.J., DeLUCIA, E.H., 2011. The greenhouse gas value of ecosystems. *Glob.*
8 *Change Biol.* 17, 425–438.
- 9 Baker, J., Ruan, X., Alcantara, A., Jones, T., Watkins, K., McDaniel, M., Frey, M., Crouse, N.,
10 Rajbhandari, R., Morehouse, J., Sanchez, J., Inglis, M., Baros, S., Penman, S., Morrison,
11 S., Budge, T., Stallcup, W., 2008. Density-dependence in urban housing unit growth: An
12 evaluation of the Pearl-Reed model for predicting housing unit stock at the census tract
13 level. *J. Econ. Soc. Meas.* 33, 155–163.
- 14 Ban, Y., Gong, P., Giri, C., 2015. Global land cover mapping using Earth observation satellite
15 data: Recent progresses and challenges. *ISPRS J. Photogramm. Remote Sens.* 103, 1–6.
- 16 Bie, C.A.J.M. de, Nguyen, T.T.H., Ali, A., Scarrott, R., Skidmore, A.K., 2012. LaHMa: a
17 landscape heterogeneity mapping method using hyper-temporal datasets. *Int. J. Geogr.*
18 *Inf. Sci.* 26, 2177–2192.
- 19 Binder, H., Tutz, G., 2007. A comparison of methods for the fitting of generalized additive
20 models. *Stat. Comput.* 18, 87–99.
- 21 Breiman, L., 2001. Random Forests. *Mach. Learn.* 45, 5–32.



- 1 Carrão, H., Gonalves, P., Caetano, M., 2010. A Nonlinear Harmonic Model for Fitting Satellite
- 2 Image Time Series: Analysis and Prediction of Land Cover Dynamics. IEEE Trans.
- 3 Geosci. Remote Sens. 48, 1919–1930.
- 4 Chaney, N.W., Sheffield, J., Villarini, G., Wood, E.F., 2014. Development of a High-Resolution
- 5 Gridded Daily Meteorological Dataset over Sub-Saharan Africa: Spatial Analysis of
- 6 Trends in Climate Extremes. J. Clim. 27, 5815–5835.
- 7 Chen, J., Chen, J., Liao, A., Cao, X., Chen, L., Chen, X., He, C., Han, G., Peng, S., Lu, M.,
- 8 Zhang, W., Tong, X., Mills, J., 2015. Global land cover mapping at 30 m resolution: A
- 9 POK-based operational approach. ISPRS J. Photogramm. Remote Sens., Global Land
- 10 Cover Mapping and Monitoring 103, 7–27.
- 11 Davin, E.L., de Noblet-Ducoudré, N., 2010. Climatic Impact of Global-Scale Deforestation:
- 12 Radiative versus Nonradiative Processes. J. Clim. 23, 97–112.
- 13 Davis, H.C., 1995. Demographic Projection Techniques for Regions and Smaller Areas: A
- 14 Primer. UBC Press.
- 15 de Beurs, K.M., Henebry, G.M., 2005. A statistical framework for the analysis of long image
- 16 time series. Int. J. Remote Sens. 26, 1551–1573.
- 17 DeFries, R.S., Field, C.B., Fung, I., Justice, C.O., Los, S., Matson, P.A., Matthews, E., Mooney,
- 18 H.A., Potter, C.S., Prentice, K., Sellers, P.J., Townshend, J.R.G., Tucker, C.J., Ustin,
- 19 S.L., Vitousek, P.M., 1995. Mapping the land surface for global atmosphere-biosphere
- 20 models: Toward continuous distributions of vegetation’s functional properties. J.
- 21 Geophys. Res. Atmospheres 100, 20867–20882.



- 1 Deichmann, U., 1996. A Review of Spatial Population Database Design and Modeling
2 (Technical Report No. 96–3). National Center for Geographic Information and Analysis,
3 Santa Barbara, CA.
- 4 Eastman, R., Sangermano, F., Ghimire, B., Zhu, H., Chen, H., Neeti, N., Cai, Y., Machado, E.A.,
5 Crema, S.C., 2009. Seasonal trend analysis of image time series. *Int. J. Remote Sens.* 30,
6 2721–2726. doi:10.1080/01431160902755338
- 7 EcoSystems Ltd, 1987. Integrated Land Use Database for Kenya. Ministry of Planning & Natural
8 Development, Nairobi, Kenya.
- 9 EcoSystems Ltd, 1983. Integrated Land Use Survey: Final Report. Lake Basin Development
10 Authority, Kisumu, Kenya.
- 11 Elzhov, T.V., Mullen, K.M., Spiess, A.N., Bolker, B., 2015. Package “minpack.lm.”
- 12 Fernández-Delgado, M., Cernadas, E., Barro, S., Amorim, D., 2014. Do We Need Hundreds of
13 Classifiers to Solve Real World Classification Problems? *J Mach Learn Res* 15, 3133–
14 3181.
- 15 Funk, C., Peterson, P., Landsfeld, M., Pedreros, D., Verdin, J., Rowland, J., Romero, B., Husak,
16 G.J., Michaelsen, J., Verdin, A., 2014. A Quasi-Global Precipitation Time Series for
17 Drought Monitoring (No. 832), U.S. Geological Survey Data Series. U.S. Geological
18 Survey, Washington, D.C.
- 19 Funk, C., Verdin, A., Michaelsen, J., Peterson, P., Pedreros, D., Husak, G., 2015. A global
20 satellite assisted precipitation climatology. *Earth Syst. Sci. Data Discuss.* 8, 401–425.
- 21 Giri, C., Pengra, B., Long, J., Loveland, T.R., 2013. Next generation of global land cover
22 characterization, mapping, and monitoring. *Int. J. Appl. Earth Obs. Geoinformation* 25,
23 30–37.



- 1 Grace, K., Husak, G., Bogle, S., 2014. Estimating agricultural production in marginal and food
2 insecure areas in Kenya using very high resolution remotely sensed imagery. *Appl.*
3 *Geogr.* 55, 257–265.
- 4 Grace, K., Husak, G.J., Harrison, L., Pedreros, D., Michaelsen, J., 2012. Using high resolution
5 satellite imagery to estimate cropped area in Guatemala and Haiti. *Appl. Geogr.* 32, 433–
6 440.
- 7 Hansen, M.C., Stehman, S.V., Potapov, P.V., 2010. Quantification of global gross forest cover
8 loss. *Proc. Natl. Acad. Sci.* 107, 8650–8655.
- 9 Hansen, M.C., Loveland, T.R., 2012. A review of large area monitoring of land cover change
10 using Landsat data. *Remote Sens. Environ., Landsat Legacy Special Issue* 122, 66–74.
- 11 Hargreaves, G.H., Samani, Z.A. 1985. Reference Crop Evapotranspiration from Temperature.
12 *Appl. Eng. Agric.* 1, 96–99.
- 13 Hastie, T.J., Tibshirani, R.J., 1990. *Generalized Additive Models*. CRC Press.
- 14 Held, I.M., Soden, B.J., 2006. Robust Responses of the Hydrological Cycle to Global Warming.
15 *J. Clim.* 19, 5686–5699.
- 16 Heistermann, M., Müller, C., Ronneberger, K., 2006. Land in sight?: Achievements, deficits and
17 potentials of continental to global scale land-use modeling. *Agric. Ecosyst. Environ.* 114,
18 141–158.
- 19 Hijmans, R.J., Phillips, S., Leathwick, J., Elith, J., 2015. Package “dismo.”
- 20 Husak, G.J., Marshall, M.T., Michaelsen, J., Pedreros, D., Funk, C., Galu, G., 2008. Crop area
21 estimation using high and medium resolution satellite imagery in areas with complex
22 topography. *J. Geophys. Res. Atmospheres* 113, D14112.



- 1 Jayne, T.S., Yamano, T., Weber, M.T., Tschirley, D., Benfica, R., Chapoto, A., Zulu, B., 2003.
2 Smallholder income and land distribution in Africa: implications for poverty reduction
3 strategies. *Food Policy* 28, 253–275.
- 4 Jayne, T.S., Muyanga, M., 2012. Land constraints in Kenya’s densely populated rural areas:
5 implications for food policy and institutional reform. *Food Secur.* 4, 399–421.
- 6 Lambin, E.F., Geist, H.J., Lepers, E., 2003. Dynamics of Land-Use and Land-Cover Change in
7 Tropical Regions. *Annu. Rev. Environ. Resour.* 28, 205–241.
- 8 Lamprey, R.H., 2013. Aerial Point Sampling (APS) Survey: Lake Basin, Machakos and Makueni,
9 Kenya, 2012-13, Nairobi, Kenya.
- 10 Lepers, E., Lambin, E.F., Janetos, A.C., DeFries, R., Achard, F., Ramankutty, N., Scholes, R.J.,
11 2005. A Synthesis of Information on Rapid Land-cover Change for the Period 1981–
12 2000. *BioScience* 55, 115–124.
- 13 Majale, M., 2008. Employment creation through participatory urban planning and slum
14 upgrading: The case of Kitale, Kenya. *Habitat Int., Labour in Urban Areas* 32, 270–282.
- 15 Makarieva, A.M., Gorshkov, V.G., Li, B.-L., 2013. Revisiting forest impact on atmospheric
16 water vapor transport and precipitation. *Theor. Appl. Climatol.* 111, 79–96.
- 17 Marshall, M.T., Husak, G.J., Michaelsen, J., Funk, C., Pedreros, D., Adoum, A., 2011. Testing a
18 high-resolution satellite interpretation technique for crop area monitoring in developing
19 countries. *Int. J. Remote Sens.* 32, 7997–8012. doi:10.1080/01431161.2010.532168
- 20 Meiyappan, P., Dalton, M., O’Neill, B.C., Jain, A.K., 2014. Spatial modeling of agricultural land
21 use change at global scale. *Ecol. Model.* 291, 152–174.



- 1 Moré, J.J., 1978. The Levenberg-Marquardt algorithm: Implementation and theory, in: Watson,
2 G.A. (Ed.), Numerical Analysis, Lecture Notes in Mathematics. Springer Berlin
3 Heidelberg, pp. 105–116.
- 4 Ngetich, K.F., Mucheru-Muna, M., Mugwe, J.N., Shisanya, C.A., Diels, J., Mugendi, D.N.,
5 2014. Length of growing season, rainfall temporal distribution, onset and cessation dates
6 in the Kenyan highlands. *Agric. For. Meteorol.* 188, 24–32.
- 7 Norton-Griffiths, M., 1988. Aerial Point Sampling for Land Use Surveys. *J. Biogeogr.* 15, 149–
8 156.
- 9 Olofsson, P., Stehman, S.V., Woodcock, C.E., Sulla-Menashe, D., Sibley, A.M., Newell, J.D.,
10 Friedl, M.A., Herold, M., 2012. A global land-cover validation data set, part I:
11 fundamental design principles. *Int. J. Remote Sens.* 33, 5768–5788.
- 12 Pielke, R.A., Pitman, A., Niyogi, D., Mahmood, R., McAlpine, C., Hossain, F., Goldewijk, K.K.,
13 Nair, U., Betts, R., Fall, S., Reichstein, M., Kabat, P., de Noblet, N., 2011. Land use/land
14 cover changes and climate: modeling analysis and observational evidence. Wiley
15 *Interdiscip. Rev. Clim. Change* 2, 828–850.
- 16 Pinzon, J.E., Tucker, C.J., 2014. A Non-Stationary 1981–2012 AVHRR NDVI3g Time Series.
17 *Remote Sens.* 6, 6929–6960.
- 18 Pitman, A.J., 2003. The evolution of, and revolution in, land surface schemes designed for
19 climate models. *Int. J. Climatol.* 23, 479–510.
- 20 Place, F., 2009. Land Tenure and Agricultural Productivity in Africa: A Comparative Analysis of
21 the Economics Literature and Recent Policy Strategies and Reforms. *World Dev.*, The
22 Limits of State-Led Land Reform 37, 1326–1336.



- 1 Platts, P.J., Omeny, P.A., Marchant, R., 2014. AFRICLIM: high-resolution climate projections
2 for ecological applications in Africa. *Afr. J. Ecol.* 53, 103–108.
- 3 Pongratz, J., Reick, C., Raddatz, T., Claussen, M., 2008. A reconstruction of global agricultural
4 areas and land cover for the last millennium. *Glob. Biogeochem. Cycles* 22, GB3018.
- 5 Pricope, N.G., Husak, G., Lopez-Carr, D., Funk, C., Michaelsen, J., 2013. The climate-
6 population nexus in the East African Horn: Emerging degradation trends in rangeland and
7 pastoral livelihood zones. *Glob. Environ. Change* 23, 1525–1541.
- 8 Rindfuss, R.R., Walsh, S.J., Turner, B.L., Fox, J., Mishra, V., 2004. Developing a science of land
9 change: Challenges and methodological issues. *Proc. Natl. Acad. Sci. U. S. A.* 101,
10 13976–13981.
- 11 Rounsevell, M.D.A., Arneth, A., Alexander, P., Brown, D.G., de Noblet-Ducoudré, N., Ellis, E.,
12 Finnigan, J., Galvin, K., Grigg, N., Harman, I., Lennox, J., Magliocca, N., Parker, D.,
13 O'Neill, B.C., Verburg, P.H., Young, O., 2014. Towards decision-based global land use
14 models for improved understanding of the Earth system. *Earth Syst Dynam* 5, 117–137.
- 15 Schaldach, R., Priess, J.A., 2008. Integrated Models of the Land System: A Review of Modelling
16 Approaches on the Regional to Global Scale. *Living Rev. Landsc. Res.* 2.
17 doi:10.12942/lrlr-2008-1
- 18 Sheffield, J., Goteti, G., Wood, E.F., 2006. Development of a 50-Year High-Resolution Global
19 Dataset of Meteorological Forcings for Land Surface Modeling. *J. Clim.* 19, 3088–3111.
- 20 Sterling, S.M., Ducharne, A., Polcher, J., 2013. The impact of global land-cover change on the
21 terrestrial water cycle. *Nat. Clim. Change* 3, 385–390.



- 1 Tian, F., Fensholt, R., Verbesselt, J., Grogan, K., Horion, S., Wang, Y., 2015. Evaluating
2 temporal consistency of long-term global NDVI datasets for trend analysis. *Remote Sens.*
3 *Environ.* 163, 326–340.
- 4 Turner, B.L., Janetos, A.C., Verbug, P.H., Murray, A.T., 2013. Land System Architecture: Using
5 Land Systems to Adapt and Mitigate Global Environmental Change. *Glob. Environ.*
6 *Change* 232395-397.
- 7 Turner, B.L., Lambin, E.F., Reenberg, A., 2007. The emergence of land change science for
8 global environmental change and sustainability. *Proc. Natl. Acad. Sci.* 104, 20666–
9 20671.
- 10 UNEP, 2008. Africa: Atlas of Our Changing Environment. UN Environment Programme,
11 Nairobi, Kenya.
- 12 Vågen, T.-G., Winowiecki, L.A., Tondoh, J.E., Desta, L.T., Gumbricht, T., 2016. Mapping of
13 soil properties and land degradation risk in Africa using MODIS reflectance. *Geoderma*
14 263, 216–225.
- 15 van Asselen, S., Verburg, P.H., 2013. Land cover change or land-use intensification: simulating
16 land system change with a global-scale land change model. *Glob. Change Biol.* 19, 3648–
17 3667.
- 18 Veldkamp, A., Fresco, L.O., 1996. CLUE-CR: An integrated multi-scale model to simulate land
19 use change scenarios in Costa Rica. *Ecol. Model.* 91, 231–248.
- 20 Verburg, P.H., Neumann, K., Nol, L., 2011. Challenges in using land use and land cover data for
21 global change studies. *Glob. Change Biol.* 17, 974–989



- 1 Verburg, P.H., Soepboer, W., Veldkamp, A., Limpiada, R., Espaldon, V., Mastura, S.S.A., 2002.
- 2 Modeling the Spatial Dynamics of Regional Land Use: The CLUE-S Model. Environ.
- 3 Manage. 30, 391–405.
- 4 Wilson, T., 2014. New Evaluations of Simple Models for Small Area Population Forecasts:
- 5 Small Area Population Forecasts. Popul. Space Place 21, 335–353.
- 6 Yu, L., Liang, L., Wang, J., Zhao, Y., Cheng, Q., Hu, L., Liu, S., Yu, L., Wang, X., Zhu, P., Li,
- 7 X., Xu, Y., Li, C., Fu, W., Li, X., Li, W., Liu, C., Cong, N., Zhang, H., Sun, F., Bi, X.,
- 8 Xin, Q., Li, D., Yan, D., Zhu, Z., Goodchild, M.F., Gong, P., 2014. Meta-discoveries
- 9 from a synthesis of satellite-based land-cover mapping research. Int. J. Remote Sens. 35,
- 10 4573–4588.

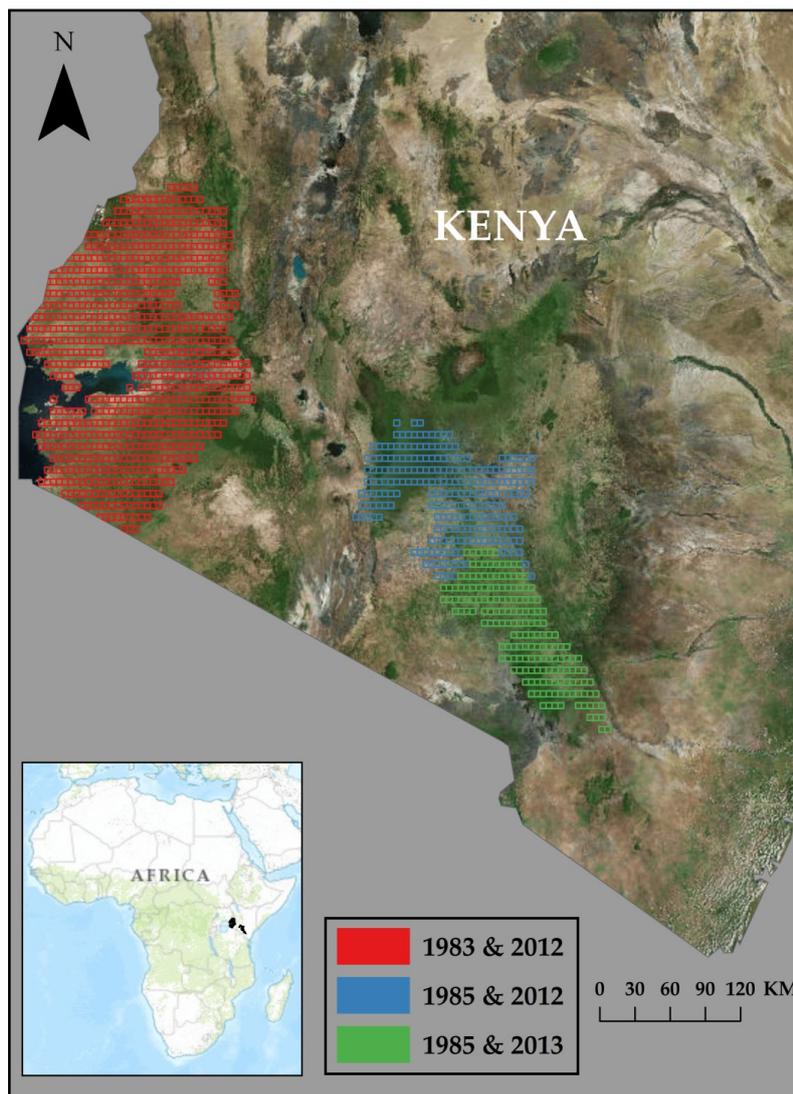


Figure 1. Study Area: 1,126 25 km² sample frames demarcating the proportion of land use/cover types estimated from aerial photo interpretation and ground surveys. Photos were taken and surveys were performed in western Kenya in 1983 and 2012, north central Kenya (Machakos area) in 1985 and 2012, and south central Kenya (Makueni area) in 1985 and 2013. Source of remote sensing image and topographic map: Environmental Systems Research Institute (ESRI).

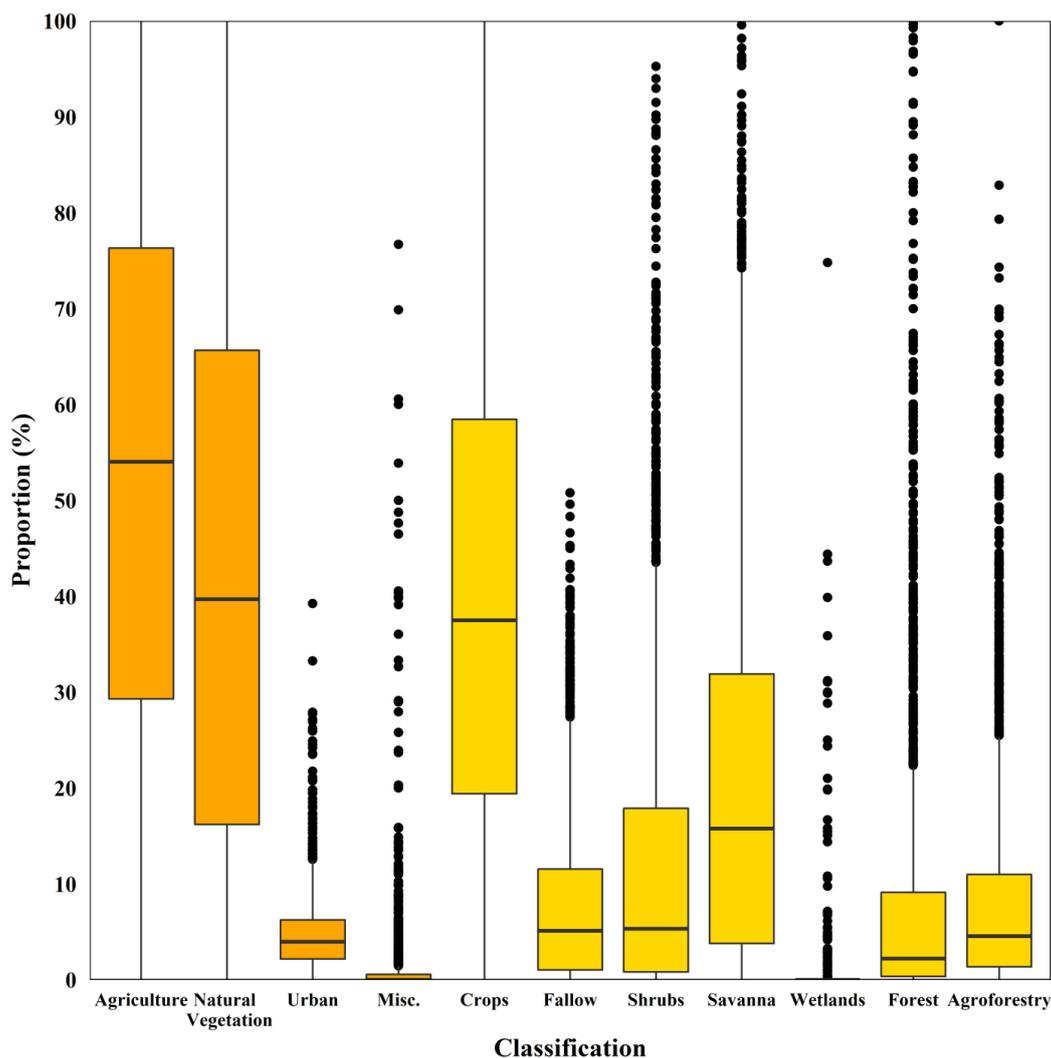


Figure 2. Boxplot of the proportion of land cover types for two levels of classification (N = 2,252). The first and second levels of classification are shaded in orange and yellow, respectively.

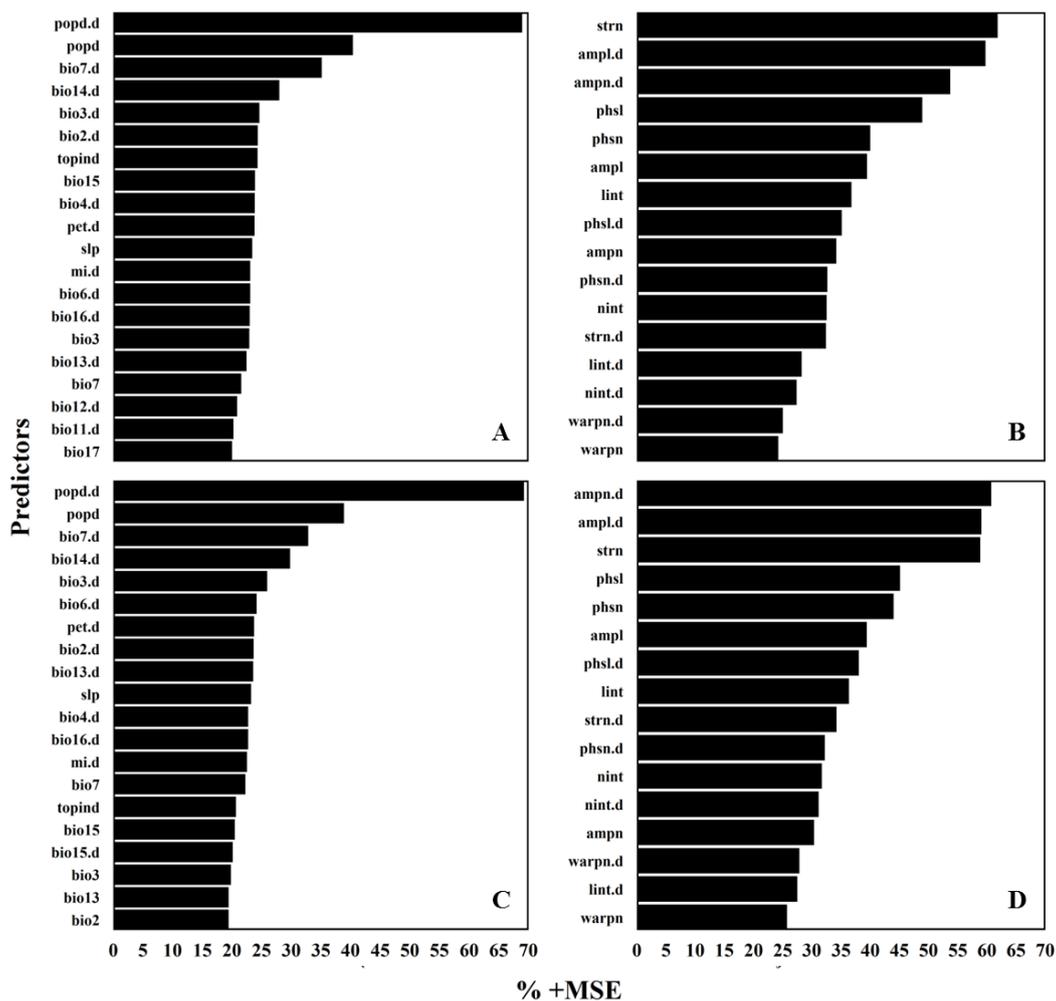


Figure 3. Percent mean squared error (MSE) increase after each of the top 20 non-remote sensing (A and C) and 16 remote sensing (B and D) predictors were omitted from the Random Forest ensemble model predicting the proportion of agriculture and natural vegetation in the calibration sample frames, respectively. The models explained 69, 49, 69, and 50% of the proportion variability.

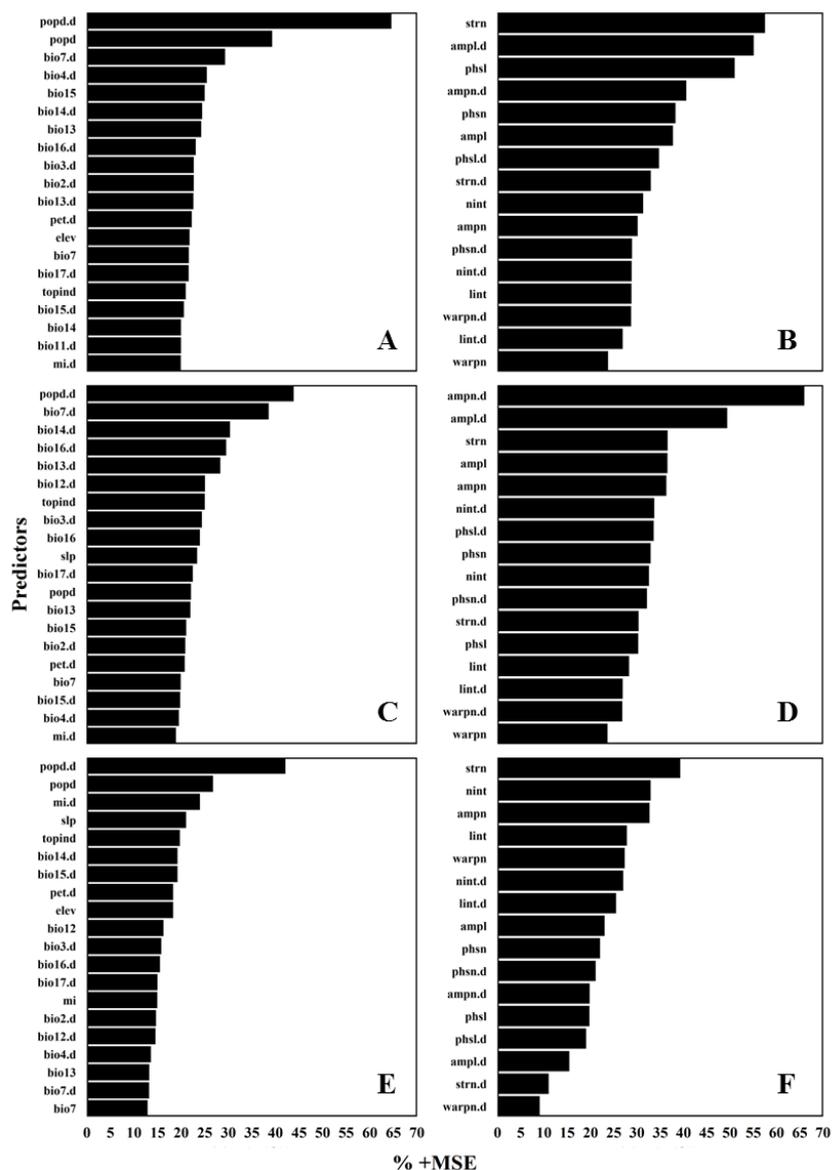


Figure 4. Percent mean squared error (MSE) increase after each of the top 20 non-remote sensing (A, C, and E) and 16 remote sensing (B, D, and F) predictors were omitted from the Random Forest ensemble model predicting the proportion of crops, savanna, and forest in the calibration sample frames, respectively. The models explained 63, 46, 62, 44, 62, and 46% of the proportion variability.

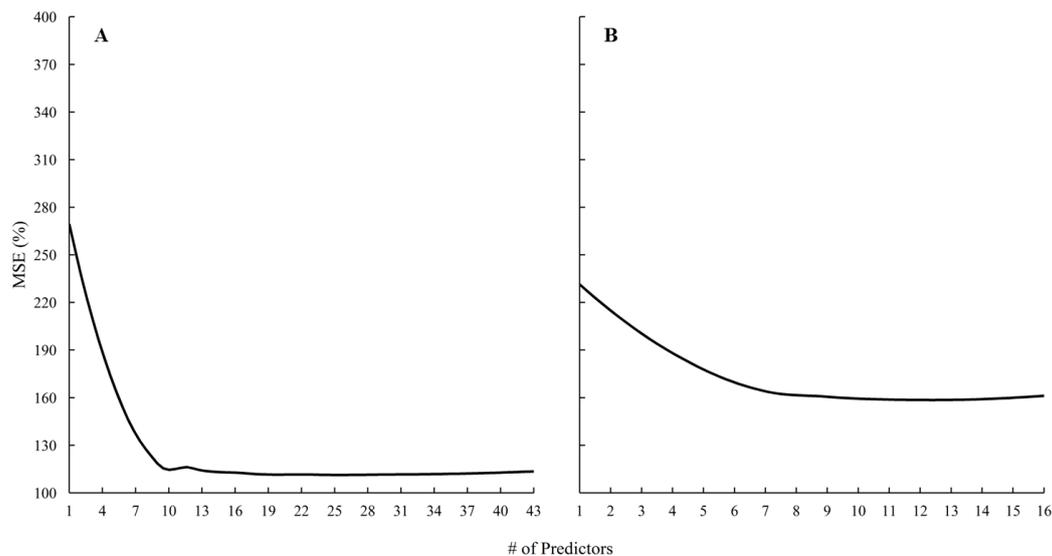


Figure 5. Curves showing the mean squared error (MSE) of the predicted proportion of forest from the Random Forest ensembles parameterized with non-remote sensing (A) and remote sensing (B) predictors. The number of predictors corresponds to the bar graphs in descending order of importance.

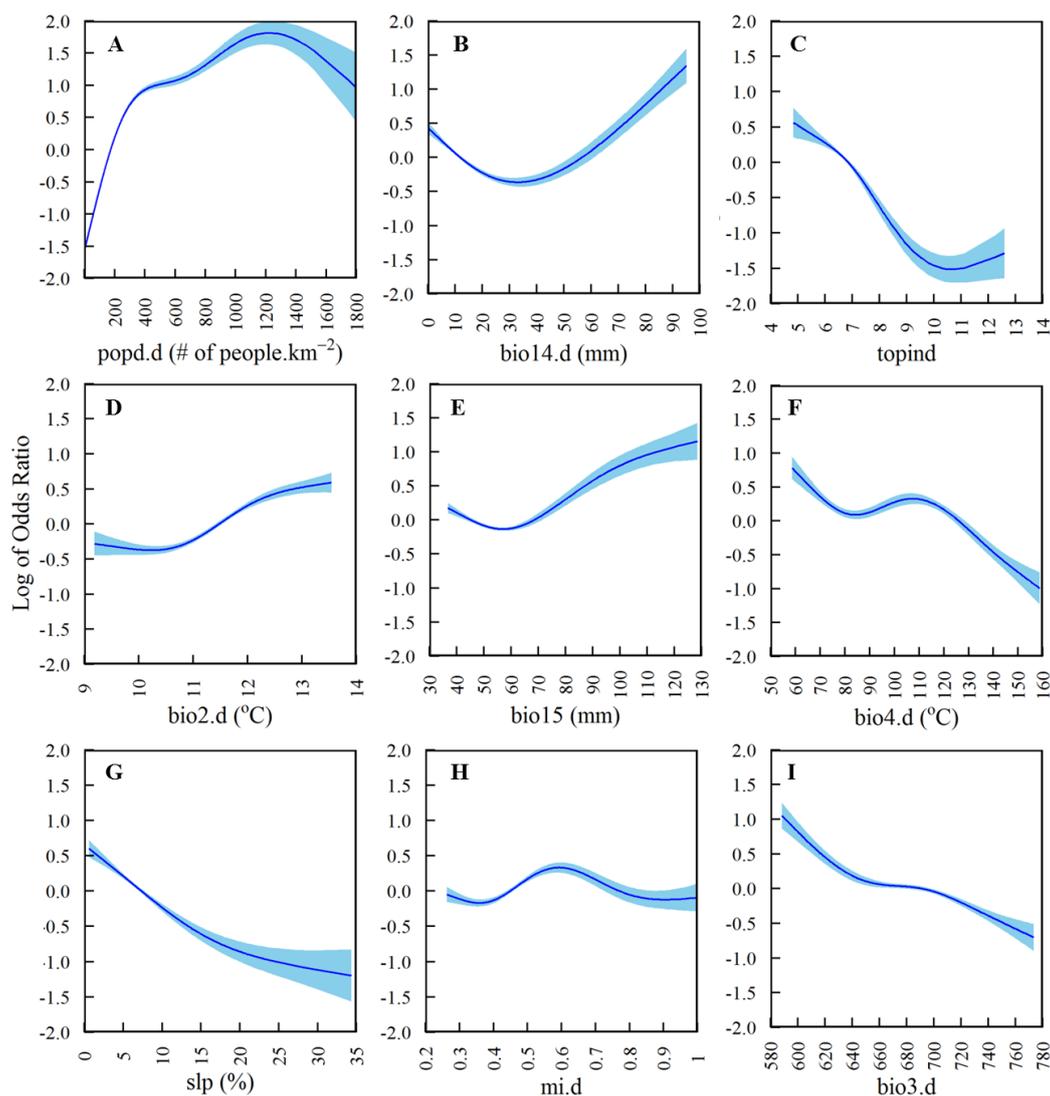


Figure 6. Partial functional plots relating the proportion (probability) of agriculture expressed as the log of odds ratio with A) population density (popd.d); B) precipitation of driest month (bio14.d); C) topographic wetness index (topind); D) mean diurnal range (bio2.d); E) precipitation seasonality (bio15); F) temperature seasonality (bio4.d); G) slope (slp); H) moisture index (mi.d); and isothermality (bio3.d). The probabilities are defined using a logistic model with cubic smoothing splines (N=1,576).

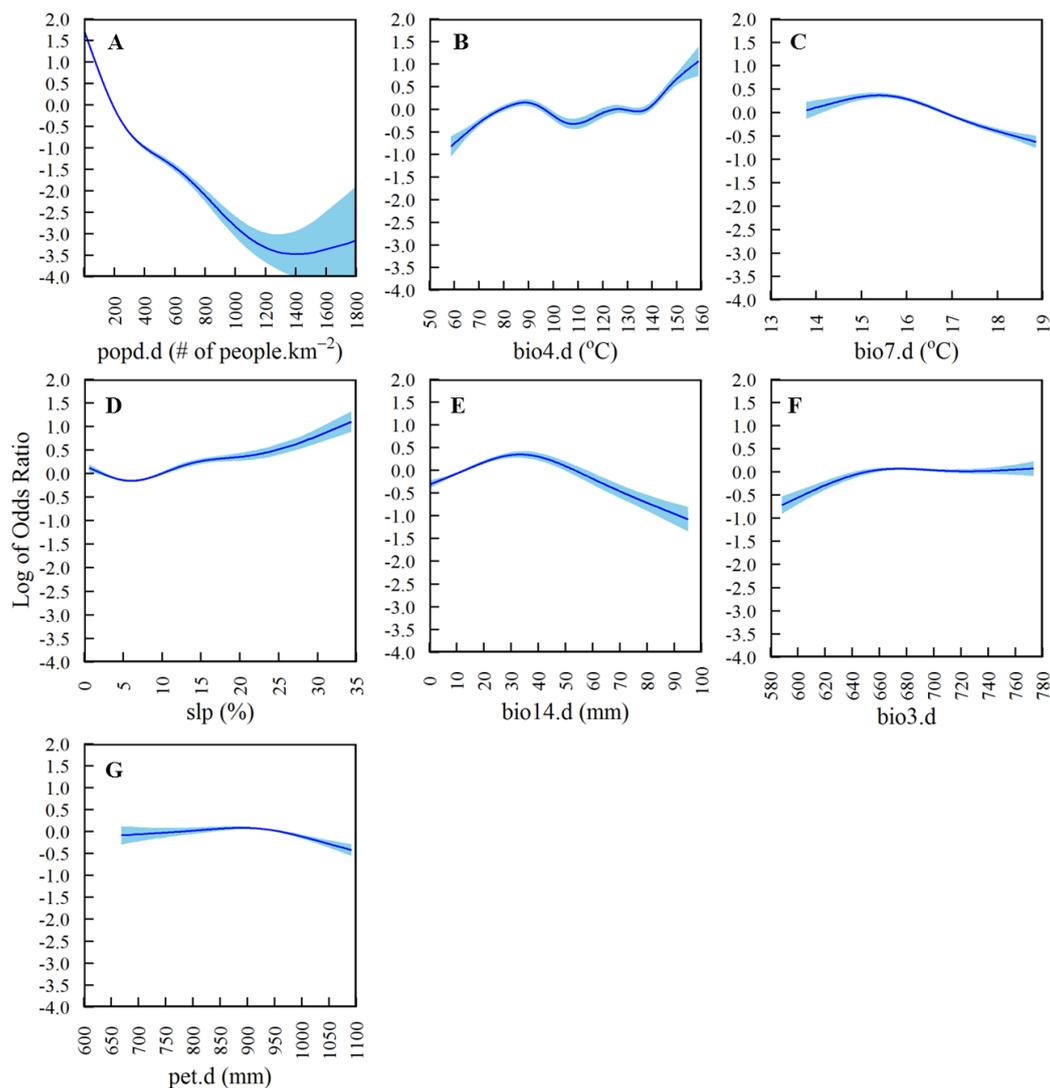


Figure 7. Partial functional plots relating the proportion (probability) of natural vegetation expressed as the log of odds ratio with A) population density (popd.d); B) temperature seasonality (bio4.d); C) temperature annual range (bio7.d); D) slope (slp); E) precipitation of the driest month (bio14.d); F) isothermality (bio3.d); and G) potential evapotranspiration (pet.d). The probabilities are defined using a logistic model with cubic smoothing splines ($N=1,576$).

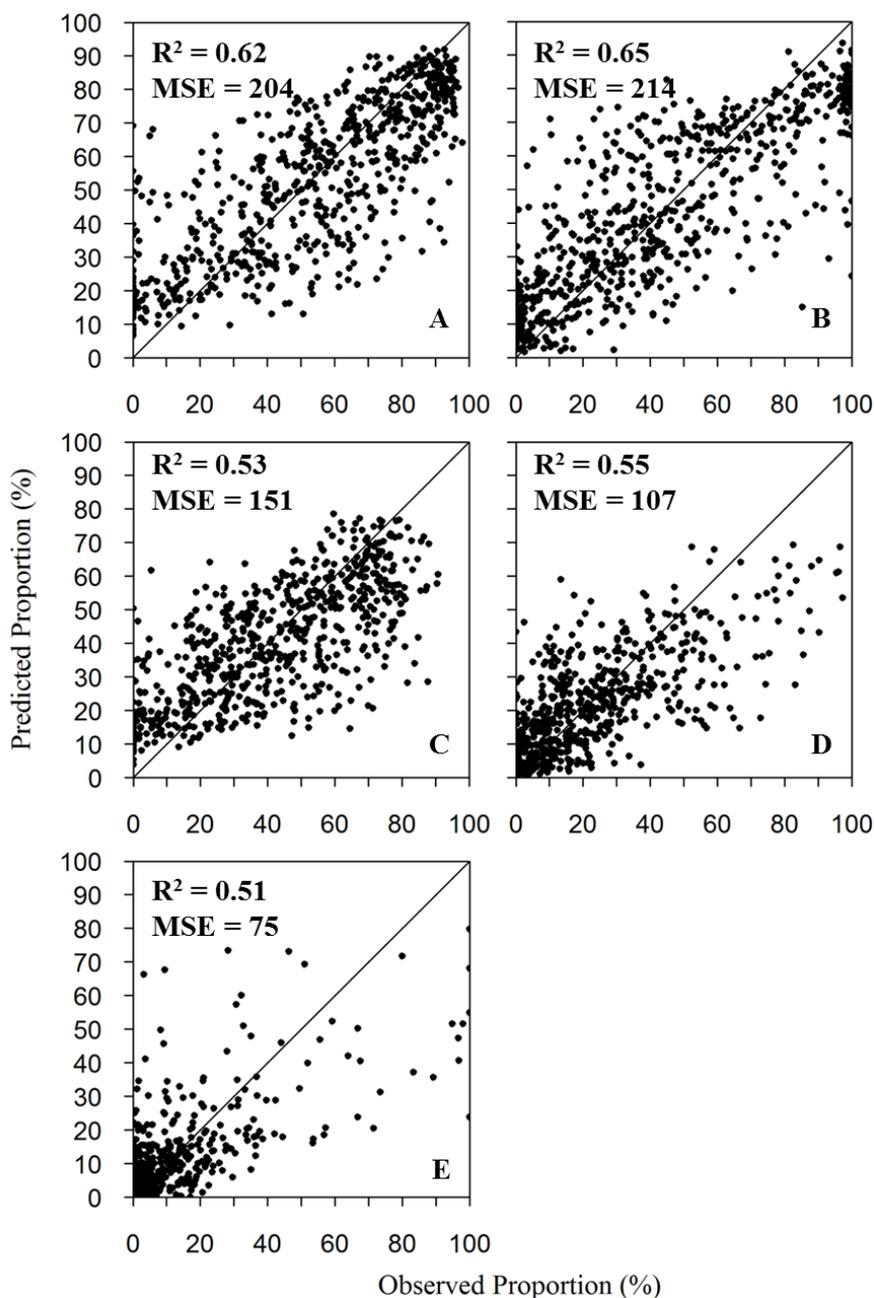


Figure 8. Predicted versus observed proportion of agriculture (A); natural vegetation (B); crops (C); savanna (D); and forest (E) for the validation subset (N=676). The 1:1 line is drawn through the origin.

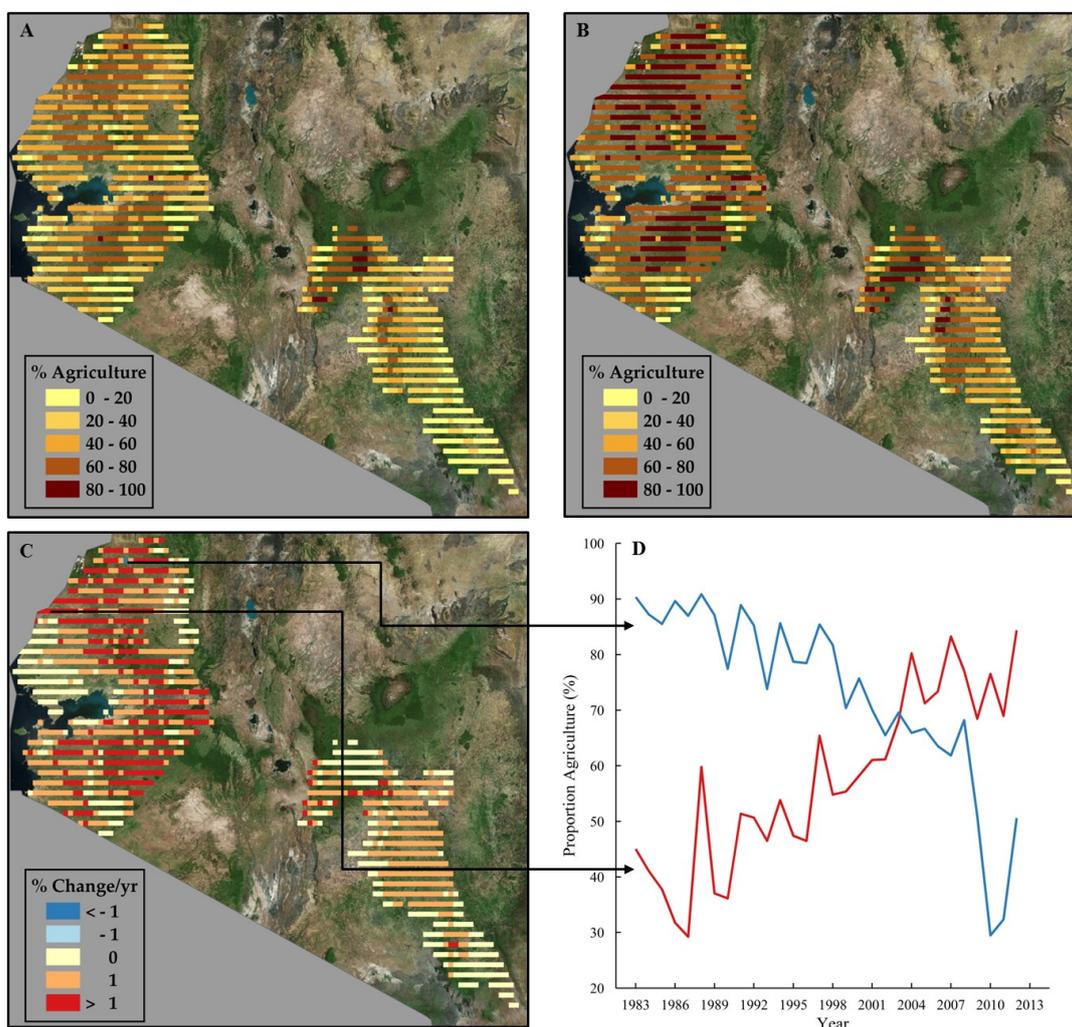


Figure 9. Simulated percent agriculture for sample frames in 1983 (A) and 2012 (B); change in agriculture per year over the 30 year (1983-2012) period (C); and time series of the strongest positive (red) and negative (blue) trend (D). Trends were determined with a Theil-Sen estimator and masked for significance using the Man-Kendall statistic at the 99.9% confidence band.



Table 1. Dates aerial sample surveys were conducted

Sample Region	First Survey	Second Survey
Lake Victoria	November 1983	October 2012
Machakos	March - May 1985	November - December 2012
Makueni	June 1985	February 2013



Table 2. Slowly-changing (long-term average/one-time value) predictors considered for LULCC estimation and their data sources. Climate, remote sensing, and population predictors were considered as annually-changing as well. Annually-changing variables are distinguished with a "d" extension.

Category	Variable	Description	Units	Source
<i>Climate</i>	bio1	Annual Mean Temperature	°C	https://www.york.ac.uk/
	bio2	Mean Diurnal Range	°C	
	bio3	Isothermality		
	bio4	Temperature Seasonality	°C	
	bio5	Maximum Temperature of Warmest Month	°C	
	bio6	Minimum Temperature of Coldest Month	°C	
	bio7	Temperature Annual Range	°C	
	bio10	Mean Temperature of Warmest Quarter	°C	
	bio11	Mean Temperature of Coldest Quarter	°C	
	bio12	Annual Precipitation	mm	
	bio13	Precipitation of Wettest Month	mm	
	bio14	Precipitation of Driest Month	mm	
	bio15	Precipitation Seasonality	mm	
	bio16	Precipitation of Wettest Quarter	mm	
	bio17	Precipitation of Driest Quarter	mm	
	mi	Moisture Index		
	pet	Potential Evapotranspiration	mm	
<i>Hydrology</i>	dtw	Depth to Groundwater	mm	http://www.bgs.ac.uk/
	gwp	Groundwater Productivity	L•s ⁻¹	
	gws	Groundwater Storage	mm	
<i>Phenological</i>	ampl	Linear Amplitude		http://ecocast.arc.nasa.gov/
	ampn	Non-linear Amplitude		
	lint	Linear Intercept (Annual Mean)		
	nint	Non-linear Intercept (Annual Mean)		
	phsl	Linear Phase		
	phsn	Non-linear Phase		
	strn	Non-linear Strength (asymmetry)		
warpn	Non-linear Warp (asymmetry)			
<i>Socioeconomic</i>	popd	Population Density	# of people•km ⁻²	http://na.unep.net/
<i>Topography</i>	asp	Aspect	°	http://www.cgiar-csi.org/
	elev	Elevation	m	
	slp	Slope	%	
	topind	Topographic Wetness Index		



Table 3. Calibration statistics of the generalized additive models used to predict the proportion of land cover (N=1,576). Predictors are significant at the 99.9% confidence band.

Land cover type	Variable ID	Part Deviance (%)	pseudo-R ²	Deviance (%)
<i>Agriculture</i>	popd.d	20.0	0.66	61.5
	bio14.d	1.9		
	topind	1.6		
	bio2.d	1.4		
	bio15	1.3		
	bio4.d	1.2		
	slp	0.9		
	mi.d	0.8		
	bio3.d	0.7		
<i>Natural Vegetation</i>	popd.d	26.2	0.66	61.4
	bio4.d	2.0		
	bio7.d	1.3		
	slp	1.2		
	bio14.d	0.6		
	bio3.d	0.5		
	pet.d	0.4		
<i>Crops</i>	popd.d	15.5	0.56	52.1
	bio2.d	3.5		
	bio15	1.8		
	bio3.d	1.7		
	bio4.d	1.4		
	pet.d	1.0		
	bio14.d	0.7		
	bio16.d	0.7		
<i>Savanna</i>	popd.d	7.0	0.56	55.7
	bio13	3.7		
	bio12.d	2.4		
	topind	2.2		
	bio16	2.2		
	bio7.d	1.6		
	bio14.d	1.6		
	bio17.d	1.4		
<i>Forest</i>	popd.d	16.4	0.57	61.2
	bio16.d	4.3		
	mi.d	2.1		



bio3.d	2.0
bio12	1.6
pet.d	1.3
elev	1.0
topind	0.7
bio14.d	0.7
slp	0.7