

1 Environmental and Ecological Statistics

2

3 Original Research

4

5 **How well does random forest analysis model deforestation and forest fragmentation in**
6 **the Brazilian Atlantic forest?**

7

8 Lisiane Zanella^{1,2,3*}, Andrew M. Folkard² George Alan Blackburn², and Luis M. T.

9 Carvalho³

10

11 ¹ Federal Institution of Rio Grande do Sul (IFRS)

12 Osório, Rio Grande do Sul, Brazil

13 e-mail: lisiane.zanella@osorio.ifrs.edu.br

14 phone: +55 51 99866 7272

15 ORCID: 0000-0002-6830-6896

16

17 ² Lancaster Environment Centre

18 Lancaster University,

19 Lancaster, Lancashire, UK

20

21 ³ Ecology and Conservation Sector, Department of Biology / Department of Forestry

22 Sciences

23 Federal University of Lavras (UFLA)

24 Lavras, Minas Gerais, Brazil

25

26 * Corresponding author

1 **Abstract**

2 We assessed the value of applying random forest analysis (RF) to relating metrics of
3 deforestation (DF) and forest fragmentation (FF) to socio-economic (S-E) and bio-
4 geophysical (BGP) factors, in the Brazilian Atlantic Forest of Minas Gerais, Brazil. A
5 vegetation-monitoring project provided land cover maps, from which we derived DF and FF
6 metrics. An ecologic-economical zoning project provided more than 300 S-E and BGP
7 factors. We used random forest analysis (RF) to identify relationships between these sets of
8 variables, and compared its performance in this task to that of a more traditional multiple
9 linear regression approach. We found that RF modelled relatively-well variance in all metrics
10 used (the rate of deforestation, the amount of forest, and the density and isolation of forest
11 patches), presenting a better performance when compared to the classical approach. RF also
12 identified geographical location and topographic factors as being most closely associated
13 with patterns of DF and FF. Both analyses found factors associated with economic
14 productivity, social institutions, accessibility and exploration to have little relationship with
15 metrics. RF was better at explaining variations in rates of deforestation, remaining forest and
16 patch patterns, than the multiple linear regression approach. We conclude that RF provides a
17 promising methodology for elucidating the relationships between land use and cover changes
18 with potential drivers.

19

20 **Keywords:** Land use and land cover change, Socio-economic and bio-geophysical factors,
21 Machine-learning technique, Stepwise Multiple Regression, Minas Gerais State, Tropical
22 forests.

23

24

25

1 **Acknowledgements**

2 The authors would like to thank the Federal University of Lavras (UFLA) for
3 providing the data. L. Zanella would like to acknowledge support from Coordenação de
4 Aperfeiçoamento de Pessoal de Nível Superior (CAPES), who provided a Ph.D. scholarship,
5 and R. Solar and T. S. de Carvalho for assistance with the statistical analysis.

6

7

1 **1 Introduction**

2 A large proportion of the Earth's surface has been transformed by anthropogenic land
3 use activities in recent centuries. Land use and land cover change (hereafter, LUCC) was
4 once considered a local environmental issue, but is becoming globally important due to its
5 increasingly widespread effects upon natural environments (Foley 2005; Lambin and Geist
6 2006). Comprehending these effects requires, in part, the understanding of relationships
7 between variations in socio-economic (hereafter, S-E) and bio-geophysical (hereafter, BGP)
8 factors associated with the LUCC with which they co-occur (Geist and Lambin 2001; Geist
9 and Lambin 2002). However, understanding these relationships is difficult because LUCC is
10 a result of complex interactions among social, economic, and environmental factors acting
11 across different scales of space and time (Geist and Lambin 2001; Geist and Lambin 2002;
12 Caldas et al. 2013). Therefore, it is necessary to design studies carefully so that inferences are
13 reliable. Unreliable conclusions can lead to distorted management recommendations,
14 resulting in missed conservation opportunities, and a waste of resources and time (Oliveira et
15 al. 2017).

16 Several studies have investigated relationships between LUCC and a wide variety of
17 S-E and environmental factors. LUCC are commonly expressed in terms of deforestation
18 rates (DF) and forest fragmentation metrics (FF). Examples of these multiscale and
19 multifactor dynamics influencing LUCC patterns are: the increasing demand for food and
20 other commodities (Aide and Grau 2004; DeFries et al. 2004; Barbier et al. 2010; Caldas et
21 al. 2013), shift in regional economies household level conditions (Perz 2004; Richards et al.
22 2008; Wright and Samaniego 2008; Gaughan et al. 2009), indirect effect of tourism (Gaughan
23 et al. 2009), globalization of markets (Hecht et al. 2006; Parés-Ramos et al. 2008), and
24 presence and effectiveness of social institutions (Hecht et al. 2006; Richards et al. 2008).

1 The studies addressing the impacts of LUCC upon tropical systems has also improved
2 significantly in recent decades (Malhi et al. 2014). Those impacts have been separated into
3 two different types: underlying (or indirect) and proximate (or immediate) causes (Geist and
4 Lambin 2002). Proximate causes are human actions that directly affect these changes, while
5 underlying causes affect these changes indirectly (Geist and Lambin 2002). The main
6 recognized proximate causes of LUCC in tropical countries are: agricultural expansion (e.g.,
7 shifting cultivation and permanent cultivation), cattle ranching, and infrastructure expansion
8 (e.g., transportation infrastructure) (Pfaff 1999; Geist and Lambin 2001; Perz et al. 2007).
9 Furthermore, LUCC is also influenced by the underlying drivers, especially demographic
10 dynamics (e.g. population growth) and economic factors (e.g. local or international demand
11 for commodities) (Geist and Lambin 2001; Caldas et al. 2013). In many regions, there is a
12 clear relationship between population change and LUCC (Geist and Lambin 2002). However,
13 other studies have shown that LUCC can be modified by socio-economic and environmental
14 factors (Geist and Lambin 2002).

15 A few studies have attempted to investigate drivers and associated factors of land use
16 and cover changes in the Brazilian Atlantic Forest. Silva et al. (2007) conducted a local scale
17 study and found an indirect influence of topographic relief on forest cover. (Teixeira et al.
18 2009) showed that proximate causes influence the dynamics of deforestation and forest re-
19 growth. They identified that losses in young secondary vegetation and forest were far from
20 rivers, on gentle slopes and near urban areas, while higher forest re-growth rates were near
21 rivers, on steep slopes and far from dirt roads. Freitas et al. (2010) analysed the effects of
22 roads, topography, and land use on forest cover dynamics and demonstrated that forest
23 dynamics were directly related to past road density, past land use (buildings and agriculture
24 expansion), and slope variation. Lira et al. (2012) described LUCC in three Atlantic Forest
25 fragmented landscapes (in São Paulo state) over time and found that LUCC deviated from a

1 random trajectory. Their results also suggested a forest transition in some Atlantic Forest
2 regions. Freitas et al. (2013) used a combination of statistical approaches – multivariate data
3 analysis (CCA), linear regression models (OLS), local spatial regression models (GWR) and
4 spatial clustering procedures (SKATER) – to investigate relationships between LUCC
5 processes and environmental and S-E factors in an Atlantic Forest region with an area of
6 ~12,000 km² in the state of Rio Grande do Sul. Their findings revealed a competitive and
7 inter-related set of LUCC processes, due to the landscape complexity. More recently, Ferreira
8 et al. (2015) investigated how forest cover and agricultural land use varied in an area of
9 Atlantic Forest in São Paulo state, emphasizing sugarcane expansion. Besides, a general trend
10 of decline followed by stabilization of forest remnants in this biome may be assumed due to
11 different deforestation rates in the Brazilian states (SOS Mata Atlântica/INPE 2014).
12 However, there are discrepancies between data sets provided by different organizations,
13 which is necessary to understand the landscape dynamics (Farinaci and Batistella 2012).

14 LUCC studies also used a range of statistical techniques. Some studies have used
15 relatively simplistic approaches, such as Mann-Whitney and Kruskal-Wallis tests (Quezada et
16 al. 2013), or correlation analyses (Beilin et al. 2014). Others have applied more robust
17 approaches, combining or comparing different methods, such as statistical redundancy
18 analyses (RDA) (Parcerisas et al. 2012); ordinary least squares regression (OLS) and
19 geographically weighted regression (GWR) (Jaimes et al. 2010; Gao and Li 2011); canonical
20 correspondence analysis (CCA), OLS, GWR and spatial clustering procedures (Freitas et al.
21 2013); stepwise multiple regression models (Gong et al. 2013). Most of these studies
22 considered a limited number of potential independent factors that had normal distributions, as
23 this is the basic requirement for using parametric techniques. Therefore, modelling
24 approaches must be further evaluated in terms of the choice of independent and metrics, as
25 well as the selection and interpretation of appropriate statistical methods. There is also a need

1 for further studies that include a large number of factors encompassing, as much as possible,
2 all aspects of the S-E and BGP context within which LUCC is taking place.

3 Despite all these improvement in our understanding of the impacts of LUCC on
4 tropical environments, there is still no optimal tool for understanding relationships between
5 deforestation/forest fragmentation and S-E or BGP factors. Random Forest analysis (RF;
6 Breiman (2001) is a variable selection technique and has great potential in this respect. RF is
7 capable of identifying complex interactive and non-linear response-predictor relationships,
8 and has excellent predictive performance (Prasad et al. 2006, Smith et al. 2011). Thus,
9 application of RF analysis to disentangle these sorts of relationships may be particularly
10 useful. RF is used widely in bioinformatics (Cutler and Stevens 2006), for land cover
11 classification (Gislason et al. 2006) and analysis of medical experiments for example, with
12 few ecological applications (Prasad et al. 2006). It has recently gained popularity in ecology
13 (Fu et al. 2010; Gilbert and Chakraborty 2011; Bonilla-Moheno et al. 2012; Ellis et al. 2012;
14 Leal et al. 2016).

15 In this study, we investigate RF regression applied to the task of identifying
16 relationships between a large set of S-E and BGP candidate independent variables (factors),
17 and metrics which quantify the current patterns of deforestation (DF) and forest
18 fragmentation (FF) of the Brazilian Atlantic Forest in the state of Minas Gerais, Brazil. This
19 study considers an unusually large set of more than 300 S-E and BGP factors. Our main
20 objective is to measure the RF ability to identify relationships with variables that describe
21 patterns of forest fragmentation, S-E/BGP and compared its results with those derived from
22 application of stepwise multiple linear regression, a classical statistical approach, to the same
23 datasets. Our hypotheses are: 1. RF is better than STEP at elucidating relationships between
24 S-E and BGP factors and FF/DF metrics. Because of RF capability of identifying complex
25 interactive and non-linear response-predictor relationships, we believe that this analysis

1 address the relationships between factor and metrics more accurately than the classical
2 approach we considered here; 2. RF and STEP identify broadly the same S-E and BGP
3 factors as being most important in explaining variation in FF/DF metrics. Based on the
4 LUCC literature, we expect that certain factors will be identified by the analyses as most
5 important, regardless of the methodological approach used, such as population and roads
6 densities, and topographic measurements (e.g. Geist and Lambin 2002; Silva et al. 2007;
7 Freitas et al. 2010).

8

9 **2 Methods**

10 ***2.1 Study area***

11 The study area is located within the state of Minas Gerais, in South-eastern Brazil and
12 comprises the 518 municipalities which fall entirely within the largest contiguous area of the
13 Atlantic Forest biome, and encompasses 34% (19,904,146 ha) of Minas Gerais (IBGE 2017,
14 Fig 1). This study site has a wide variability across the municipalities in the magnitude of
15 DF/FF metrics and in the S-E/BGP factor values.

16 The study region characterized by rolling hills which rise from 200 m to a medium
17 altitude of 1600 m. It is a very rugged area with a large proportion of highlands as well as
18 plateaus and plains. There are several climates types linked to the different relieves: warmer
19 climate in the north and cooler in the south. The distance from the ocean also has a climatic
20 effect (maritime vs. inland climate, etc) upon the study area. The region is, on average,
21 relatively sparsely populated, with a tendency in higher concentrations of populations
22 towards the south, which has smallest municipality areas too. The south part of the study area
23 is also relatively richer and developed when compared to the other part of the study area and
24 to the Brazilian average. The main industries and sources of employment are the bovine cattle
25 herd - which corresponds to 10% of the Brazilian total -, coffee production and the extraction

1 of iron ore. The location of Belo Horizonte, the largest city and the capital of Minas Gerais,
2 also plays an important role in the establishment of many industries, especially automobile
3 and steel mill industries in its vicinities. This makes Minas Gerais the second largest
4 automotive and metal foundry hub in Brazil. All these information on the study area patterns
5 and more can be found in the ecologic-economical zoning of Minas Gerais, ZEE-MG
6 (Scolforo et al. 2008).

7

8 #Fig 1 approximately here

9

10 **2.2 Variable selection**

11 This work used large datasets provided by two broader-scale projects carried out in
12 Minas Gerais State, Brazil. The DF and FF metrics were derived from the vegetation
13 monitoring system dataset (Scolforo and Carvalho 2006; Carvalho and Scolforo 2008,
14 Carvalho and Scolforo - *unpublished data*), which comprises land cover maps from 2003 to
15 2011.

16 A deforestation metric, the growth rate of deforestation (GRD, percentage) was
17 calculated for each municipality using digital change detection applied to Landsat images
18 from the vegetation monitoring system dataset (Scolforo and Carvalho 2006; Carvalho and
19 Scolforo 2008, Carvalho and Scolforo - *unpublished data*). GRD was normalized by the
20 amount of remaining forest area within each municipality.

21 To quantify forest fragmentation, we used the 2011 land cover map from the
22 vegetation monitoring system dataset (Scolforo and Carvalho 2006; Carvalho and Scolforo
23 2008, Carvalho and Scolforo - *unpublished data*). A set of 225 landscape metrics from class
24 and landscape levels from all of the different categories available in FragStats 4.0 (McGarigal
25 et al. 2012) were calculated for each of the 518 municipalities considering the forest cover

1 configuration in 2011. These were then passed through a three-stage filtering process to
2 provide a tractable set of metrics for use in our analysis of statistical approaches. Firstly,
3 noting that metrics in datasets such as this can be highly correlated (Riitters et al. 1995), we
4 selected a subset of uncorrelated metrics based on *Pearson* correlation analyses conducted
5 using the Pairs-panel analyses in R. We discarded those metrics which were strongly
6 correlated (defined for these purposes as having correlation coefficients for which $p \leq 0.01$)
7 with selected variables, and therefore deemed to be redundant. When two or more variables
8 were significantly correlated, the selection criteria to choose one of them were mathematical
9 simplicity and an intuitive judgment of their explanatory power in terms of ecological
10 meaning. Secondly, we chose metrics from the remaining subset that were commonly used in
11 literature (those which were repeatedly found in the papers consulted) found via a search on
12 the Web of Knowledge website (<http://wok.mimas.ac.uk/>). The search was carried out from
13 2011 to June 2013, using the key-words "landscape metrics" and/or "landscape indices". This
14 search yielded 48 papers, of which four were found, on inspection, to be out of scope, and we
15 had no access to another five. The papers consulted in the review can be seen in the
16 Supplementary material (List S1 – ESM1). Finally, we verified the normality of the residuals
17 from linear models (see the section *Stepwise multiple linear regression* for more details) and
18 those metrics which had non-normally distributed residuals were discarded to enable
19 comparative analysis of the random forest method with classical, parametric multiple
20 regression, which requires normally distributed variables most of the times. At the end of the
21 three-stage filtering process, three landscape metrics representing forest fragmentation at
22 municipality scale were selected: the total remaining forest (CA), a measure of forest cover;
23 the mean Euclidean nearest-neighbour distance (ENN), a measure of patch's isolation from
24 each other; and the patch density (PD), a measure of forest spatial structure (Table 1).

1 The S-E and BGP factors were derived from the ecologic-economical zoning of
2 Minas Gerais, ZEE-MG (Scolforo et al. 2008). Almost all available factors were derived
3 within political administrative units at the scale of municipalities, the smallest administrative
4 units in Brazil. To avoid bias, we chose to use only the metrics that would allow us to analyse
5 them at the municipality scale.

6 S-E and BGP factors were obtained from the ZEE-MG database, which collates data
7 from different national agencies. The years for which these variables were collected were
8 limited by the availability of information from national agencies, and ranged from 2003 to
9 2006. Based on data availability, socio-economic factors from four categories – production,
10 exploration, human and institutional – were used. Variables from further four categories of
11 BGP factors – topography, distance, accessibility, and geographical location – were also
12 selected. This gave an initial list of more than 300 candidate independent factors.

13 Descriptions of how these variables were calculated can be found in Scolforo et al. 2008.

14 From this list, a tractable sub-set of factors was derived using the first step from the filtering
15 process described above for the FF metrics. As a result, a total of 34 S-E and BGP factors
16 were selected as factors for use in our comparative analysis of statistical approaches (see
17 Table S2, in the supplementary material, for a complete description of all factors).

18

19 ***2.3 Random forest analysis (RF)***

20 Random forest analysis is a machine-learning technique that may be used for
21 predictive modelling of multiple outcomes from large input datasets. In short, RF uses an
22 ensemble of decision trees with binary divisions, each capable of producing an outcome when
23 presented with a set of input values (Cutler et al. 2007). For regression modelling problems
24 the tree response is an estimate of dependent (outcome) variable values derived from the
25 given values of a set of independent (input) variables. RF uses a regression tree approach

1 (also known as "CART"; Breiman et al. 1984), to build a number of decision tree models
2 from randomly selected subsets of training samples and factors (Cutler et al. 2007). Model
3 fitness is examined using validation data that is not in the training sub-sample; hence, cross-
4 validation with external data is not necessary. The validation sample is also used to calculate
5 measures of variable relative importance (Ellis et al. 2012). The outcomes from all of the
6 trees are then averaged, which provides predictive accuracy and low bias (Breiman 2001).

7 We used the R "extendedForest" library provided by the Gradient Forest project
8 (Smith et al. 2011; Ellis et al. 2012) to carry out RF analysis. This package was developed for
9 use in ecological studies of species distributions. It integrates results from RF analyses for a
10 number of individual species distributions into results that enable prediction of multiple
11 species distributions (Smith et al. 2011; Ellis et al. 2012). In addition, it is able to analyse
12 large numbers of potential factors and to reduce bias when predictors are correlated (Smith et
13 al. 2011). In our study, we extended the application of extendedForest by using the DF and
14 FF metrics described above (i.e. GDR, ENN, CA and PD) in place of the species distributions
15 used in the application for which it was originally developed. We build partial dependence
16 plots using the variable relative importance values. Models were fitted with 10,000 trees. In
17 each split, we used one-third of the factors randomly sampled as independent candidates. We
18 excluded from final models the variables with negative relative importance values, which do
19 not contribute to the overall explanation. In order to test our first hypothesis, we also
20 calculated the R² in RF approach to compare it with outcomes from the stepwise multiple
21 linear regression.

22

23 ***2.4 Stepwise Multiple Linear Regression***

24 From a wide range of possible approaches, we selected stepwise multiple linear
25 regression (hereafter, STEP) as a comparator method against which to assess the performance

1 of RF. This type of technique is arguably the most common approach to data-based
2 prediction and simulation tasks (Whittingham et al. 2006). For situations in which the number
3 of variables is high, as is the case here, it is appropriate to incorporate into the modelling
4 process a method for selecting only those factors that contribute most strongly to the
5 predictive model delivered. The STEP approach to multiple regression is a routine technique
6 for achieving this (see, for example, Efroymson 1960; Hocking 1976; Furundzic 1998).
7 Despite having a number of weaknesses, notably bias in parameter estimation,
8 inconsistencies among model selection algorithms, and an inappropriate focus on a single
9 best model (Burnham and Anderson 2002; Kadane and Lazar 2004; Whittingham et al.
10 2006), it is used widely within ecology and landscape studies (Whittingham et al. 2006).

11 The stepwise method combines forward selection and backward elimination
12 procedures (Venables and Ripley 2002; James et al. 2013). It proceeded by first setting up an
13 initial model incorporating a subset of the candidate independent variables (factors). Then,
14 this model was iteratively altered by adding significant factors and/or removing insignificant
15 ones, in a process called the stepping procedure. A variable that enters at an early stage may
16 become superfluous at later stages because of its relationship with other factors subsequently
17 added to the model (Kleinbaum et al. 1998). To check this possibility, at each step a partial F
18 test is carried out for each factors currently in the model, regardless of the stage at which it
19 was entered. The whole process is repeated until no more factors can be added or removed,
20 which means that the model is optimized, or when a specified maximum number of steps is
21 reached. Many statistical methods are available to test the stability and validity of the final
22 regression model. We used the adjusted square of the correlation coefficient (adjusted R^2) and
23 the AIC (Akaike Information Criteria) to assess our final model. The AIC was also used to
24 calculate relative variable importance. Implementation was based on the dredge function for
25 automated model selection, which is available as the R “MuMIn” package (Barton 2014). It

1 calculates AIC values for models with all possible combinations of factors and ranks the
2 models based on the calculated values. MuMin is also highly demanding in terms of
3 computational time and resource requirements. We determined the relative importance of
4 each independent variable selected in the models from STEP approach based on AIC weights
5 (importance function in MuMIn; Burnham and Anderson 2002). The relative importance
6 values were converted to percentages for comparison with the equivalent outcomes from RF.

7

8 ***2.5 Final models***

9 We used specific acronyms for the models we have tested to make it easier for readers
10 to understand them. For this, we use the acronyms of each of the metrics tested, which reflect
11 deforestation (DF): GRD; and forest fragmentation (FF): CA, ENN and PD and we add the
12 acronym of the two analysis approaches that we used: RF and STEP. The results were four
13 models selected using RF approach and four other using STEP approach, respectively: the
14 growth rate of deforestation – RF-GRD and STEP-GRD; the total remaining forest – RF-CA
15 and STEP-CA; the mean Euclidean nearest-neighbour distance – RF-ENN and STEP-ENN;
16 and the patch density – RF-PD and STEP-PD models.

17

18 **3 Results**

19 ***3.1 Random forest analysis***

20 The RF analysis provides evidence of the effect SE and BGP factors (see Table S2 in
21 the supplementary material ESM2) on the metrics, explaining high amounts of the observed
22 variance (up to 99%) of some of them, and lower amounts of the observed variance of others
23 (less than ~ 40%) (Fig 2 - see also Table S3 in the supplementary material ESM3). In the
24 latter cases, the outcomes imply that there is restricted explanatory power in the factors, and
25 that variability in some of the models across the municipalities is not explained by the factors

1 considered here. The relative importance of each factor was quantified as its partial
2 contribution to explaining the variability of each of the four metrics tested by both statistical
3 approaches, expressed as a percentage. Although, these values are not quantitatively
4 comparable between the metrics, they allow us to rank the factors in terms of their relative
5 importance in each metric model.

6

7 #Fig 2 approximately here

8

9 Of the four models using RF approach, RF-GRD performed best, with a very high
10 value (99.40% - Fig 2) of its variance explained by the factors. Distance variables (longitude
11 and the minimum distance of forest patches to the nearest reservoir and the nearest protected
12 area) and geographical location were the most important factors in this respect. Among the
13 many factors selected in GRD model selected by RF, those related to topography and crop
14 production were also relatively important. Longitude (POINT_X) explained a greatest part of
15 the variance in RF-GRD model (Fig 3.a)

16

17 #Fig 3 approximately here

18

19 The selected patch density model (RF-PD), had the second highest amount of its
20 variation explained (61.52%, Fig 2). A large number of factors were identified as having
21 some role in explaining RF-PD variations between municipalities; those with the highest
22 importance were associated with the road network or were topographic. Roads density was
23 the factor which most explained the variance in this model (Fig 3.b).

24

25 The selected models of total remaining forest (RF-CA) and of the mean Euclidean
nearest-neighbour distance between forest patches (RF-ENN) also had relatively-high

1 amounts of their variation explained (40.67 and 39.38%, respectively, Fig 2). The factors
2 with the highest importance for predicting these models were the mean slope of each
3 municipality (Fig 3.c) for the selected RF-CA model and the mean altitude of each
4 municipality (Fig 3.d) for the selected RF-ENN model. Other topographic factors (the mean
5 altitude across each municipality for RF-CA, and the mean slope across each whole
6 municipality, and the mean slope within deforested areas, for RF-ENN) were also relatively
7 important, as were geographical location, distances to the nearest protected area and nearest
8 steel mill, and longitude.

9 Overall, factors from the geographical location, distance, topography, institutional and
10 accessibility categories appeared among the most important factors in all the four selected
11 models from RF approach, namely: the latitude of municipalities; the minimum distance from
12 forest patches to the nearest steel mill and the longitude of municipalities; mean slope, mean
13 slope within deforested areas and mean altitude; the amount of protected area in each
14 municipality; and the density of roads.

15

16 ***3.2 Comparisons of RF with STEP***

17 Outcomes from the STEP approach are shown alongside those for RF, in as
18 comparable a form as possible (Fig 2). Note that, although "percentage importance" values
19 are quoted for models from both analysis approaches, these values are not quantitatively
20 comparable between these two methods' outcomes or between different metrics addressed in
21 models. Rather, these values allow us to rank the factors in terms of their relative importance
22 for explaining the variability of each model. The percentages of variance explained by the
23 two analysis approaches are, however, comparable. Both approaches provided evidence of
24 relevant relationships, but models from RF approach surpassed the capacity of the classical

1 approach in explain models' variance. However, the results are mixed in terms of the factors
2 selected as being most important by each approach.

3 The selected STEP-CA model performed best of all models from STEP approach. It
4 explained an amount (39.80% c.f. 40.67% for RF-CA) of CA variation between
5 municipalities similar to that explained by RF. There was also a strong similarity between the
6 most important factors selected by the models from both approaches, since all of the factors
7 selected by STEP were also selected by RF, except soil types and employability. The mean
8 slope was the most important factor explaining the selected models from both approaches.
9 Other important factors were latitude, longitude and mean altitude. The amount of protected
10 area in each municipality and the number of rural family farms were also important in STEP-
11 CA.

12 STEP-ENN had the second highest value of ENN explained variance f (30.91% by
13 STEP-ENN, 39.38% by RF-ENN). Factors were less similar between ENN models than in
14 the CA models. While the mean altitude was the most important factor found by RF-ENN,
15 four factors were important in the STEP-ENN selected model, namely: the mean slope, soil
16 type, density of roads and latitude.

17 The selected PD model from STEP approach (STEP-PD) also had a relatively high
18 amount of its variance explained compared to the other models from STEP, but much less
19 than the selected RF-PD model (29.40% c.f. 61.52% for RF-PD). Some of the factors were
20 found in the selected models from both approaches. However, only one of the most important
21 factors appeared in both of these models: the mean slope of deforestation patches, a
22 topographic factor. The density of roads was the factor identified as being most important by
23 RF-PD, while a similar factor, the minimum distance to the nearest road had the highest
24 importance in STEP-PD. Another topographic factor important in the STEP-PD was the

1 minimum mean slope within each municipality, while in RF-PD the mean altitude, and
2 latitude were also important.

3 There was a strong contrast between the amounts of variance explained for the growth
4 rate of deforestation by STEP (17.36%) and RF (99.4%) approaches. In STEP-GRD, the
5 minimum distances to the nearest protected area and nearest steel mill were the most
6 important factor explaining GRD variance, followed by the mean slope and the amount of
7 protected area. In RF-GRD, the longitude and, secondarily, the latitude and minimum
8 distances to the nearest steel mill and nearest reservoir were also important.

9

10 **4 Discussion**

11 *4.1 Random Forest analysis*

12 In the RF approach' outcomes, we observed that there are some strong relationships
13 between the S-E and BGP factors and DF and FF metrics. RF performed best for the growth
14 rate of deforestation (RF-GRD) and secondarily for patch density (RF-PD) selected models,
15 explaining around 99% and 60% of their variances, respectively – high values for ecological
16 studies. It also performed relatively well for the total remaining forest (RF-CA) and patch
17 isolation mean Euclidean nearest neighbour distance (RF-ENN) selected models, explaining
18 40.67% and 39.38% of their variances, respectively. In terms of model performance, this may
19 suggest that the random forest approach is good at identifying parameters that describe some
20 macro-scale factors (rate of deforestation and the overall remaining forest) and the
21 distribution of patches within a landscape (their density and mean isolation from each other.
22 Alternatively, these results could be interpreted as indicating that the rate of deforestation,
23 remaining forest and patch-distribution scale variables (GRD, PD, CA and ENN) are closely
24 linked to the factors we have considered here. In other words, RF is particularly good at
25 identifying links for the types of parameters we analyse, since it performs better providing a

1 higher amount of metrics variance explanation. It is important to note that, even using a very
2 large dataset comprising many factors, much of the variance in some of the four metrics was
3 not accounted by our selected models. In addition, the question of whether it is primarily the
4 nature of the model or the nature of the factors that has led to this finding is not answerable
5 by this first application of RF to this type of data, and remains to be addressed by further
6 investigation.

7 Turning now to consideration of the factors, we found that some of them were
8 particularly strongly related to some of the metrics, for example longitude (which explained
9 20.7% of GRD variance), road density (which explained 20.4% of PD), and mean altitude
10 (which explained 18.5% of ENN). However, neither the nature of, nor the reason (i.e.
11 whether they are causatively-linked or simply co-vary) for these links are elucidated by RF.
12 Despite these cases of strong individual-variable links, no single independent variable was
13 found to be related to all of the metrics. Geist and Lambin (2002), who investigated the
14 causes of deforestation of tropical forests, also did not find a single important factor. They
15 concluded that forest loss is due to a combination of factors that vary with historical and
16 geographical context. We conclude from the present study that we can expect the same for
17 forest fragmentation.

18 At the level of independent variable categories and considering only the three
19 variables in each model which made the strongest contributions the explain metrics variance,
20 we found that those from the Geographical location, Topography, Distance and Accessibility
21 categories contributed most to explaining variance in the RF outcomes. On the other hand,
22 variables from the Exploration, Institutional, and Productivity categories made hardly any
23 contribution. Additionally, we found that factors from the Geographical location and
24 Topography categories made up the majority of the most-important independent variable
25 explaining each dependent variable in the models from RF approach. This suggests that the

1 physical environment is more important for determining variations in DF and FF metrics
2 between municipalities, than social or economic issues. Other studies conducted in the
3 Atlantic Forest agree with our results, showing that physical environment factors play a
4 significant role on deforestation and forest fragmentation (Silva et al. 2007; Teixeira et al.
5 2009; Freitas et al. 2010). In other countries of Latin America, a similar pattern can be also
6 observed, with physical environment being more important than socioeconomic or
7 demographic factors to explain land-cover change (Bonilla-Moheno et al. 2012; Redo et al.
8 2012). In addition, specifically in our case, geographical location is important considering the
9 discrepancies between the north and south parts of the study area, mainly in terms of
10 development, what also could work as a proxy of some socioeconomic and demographic
11 factors. However, these findings do not exclude the contribution of socioeconomic or
12 demographic factors upon deforestation and forest fragmentation, since they might be
13 indirectly linked to the physical environment factors. For example, deforestation is more
14 likely to be located in lower and less steep terrain, where transport and mechanical
15 agriculture are easier (Apan and Peterson 1998). They are more likely to have occurred in
16 sites more suitable for agriculture (Flamenco-Sandoval et al. 2007; Killeen et al. 2007;
17 Fearnside 2016). This finding has important implications for management policies aimed at
18 conserving the Atlantic forest, and possibly other biomes that are fragmenting under
19 anthropogenic pressures, although it requires further evidence to be confirmed. This points
20 out the importance of valuing biodiversity in impacted sites (lower and less steep terrain)
21 when selecting areas for conservation, for example (Margules and Pressey 2000; Metzger and
22 Casatti 2006). Also, although this ordering of importance of the different types of factors is
23 quite coherent across the RF approach' outcomes, the question remains as to whether it is
24 "true". Claims to this effect are supported by noting that factors that random forest-type

1 methods have identified as most important for classification have been found to coincide with
2 ecological expectations in the literature (Cutler et al. 2007; Wei et al. 2010; Ellis et al. 2012).

3

4 ***4.2 Comparisons of RF with STEP***

5 Like RF, the STEP approach found some strong relationships between the S-E/BGP
6 factors and DF/FF metrics. Unlike RF, STEP selected models found the most explained-
7 variance and strongest relationships for the amount of forest, followed by the isolation of
8 forest patches. Unlike RF, however, there was less difference in the performances of models
9 from STEP approach: while the explained variances from RF ranged from ~40% to 99%,
10 STEP explained between ~18 and 40% of the variance of all four metrics, confirming our
11 first hypothesis, that RF addresses the relationships between factor and metrics more
12 accurately than STEP approach.

13 Contrary to our second hypothesis, there was more disagreement than agreement,
14 overall, in terms of the selection and importance of the factors between the two approaches.
15 A low number of factors was selected as important and shared by them. Considering the
16 categories of factors, both approaches found that factors from the Topography category were
17 of higher importance in all selected models, while the Geographical location was more
18 important in the selected models from RF than from STEP approach. Variables from the
19 Distances and Accessibility categories were of intermediate importance, and variables from
20 the Exploration, Institutional, and Production categories were of little importance. In the
21 selected models from STEP approach, we found that the most-important independent variable
22 explaining each dependent variable model also belonged to the Distances and Topographic
23 categories.

24 The most important factors of selected models in RF approach were subtly different
25 than those selected in STEP approach. Considering the selected rate of deforestation model

1 from RF, the most important factor influencing it is longitude of municipalities, which
2 represents a measure of the distance from the ocean (climate) and also to socioeconomic
3 longitudinal gradient. We expected that deforestation increases towards a socioeconomic
4 gradient, which may reflect a higher degree of developed, and consequently, higher
5 exploration of natural resources, for example. On the other hand, the most important factors
6 in the selected model from STEP were the minimum distance to the protected area. In a
7 similar way, we expected that deforestation decreases when forest patches are closer to
8 natural reserves. The smaller the distance, the closer the forest patches are to a natural
9 reserve. This may mean that there is a greater amount of forest in the municipalities where
10 the forest patches are closer to the natural reserves, whereas in those municipalities where the
11 reserves are more distant, there is possibly a smaller amount of forest, and therefore,
12 deforestation rate is also smaller. Although different, these two factors may be ecologically
13 linked to deforestation rates.

14 Turning to isolation of forest patches, two different factors from the Topography
15 category appeared as most important factors in the selected models from RF and STEP,
16 respectively, the mean altitude of each municipality and the mean slope across each whole
17 municipality. Although different measurements, these factors are related to the relief of the
18 study area, that plays an important role influencing deforestation (Silva et al. 2007) in The
19 Atlantic Forest Biome. Also, due to an intense exploration in the last 500 year, the Atlantic
20 forest remnants are currently restrict to the higher elevations and steeper reliefs (Dean 1996;
21 Oliveira-Filho and Fontes 2000; Ribeiro et al. 2009; Kauano et al. 2012). The most important
22 factor was the same for the amount of forest in both selected models from RF and STEP
23 approaches: mean slope; also related to the study area relief.

24 The density of forest patches was mostly affected by two similar factors: the density
25 of roads in the selected model from RF; and the minimum distance to the nearest road in the

1 selected model from STEP. These findings are consistent, since roads serve as fragmenting
2 features (Forman and Alexander 1998; Butler et al. 2004), subdividing forests, increasing the
3 number of forest patches, and reducing forest connectance. Roads have few positive, neutral
4 and numerous negative environmental impacts. Positive impacts include increasing
5 accessibility (Leinbach 1995), which can also be negative since this facilitates deforestation
6 (Laurance et al. 2001). Negative impacts include habitat loss, degradation, and fragmentation,
7 direct wildlife mortality, and road avoidance behaviours by wildlife (Forman and Alexander
8 1998). Therefore, density of roads plays an effective role in forest fragmentation, and the
9 minimum distance to the nearest road also reflects this role.

10 Notwithstanding a few similarities between the outcomes of the two modelling
11 approaches, differences between them are strongly evident. However, the reasons for these
12 differences are not clear from our results, and require further investigation. Nonetheless, in
13 theory, one would expect the RF approach' outcomes to identify more reliably than STEP the
14 factors that have greatest influence over models. This expectation arises from the greater
15 robustness of random-forest type methods compared to traditional regression approaches.
16 Unlike traditional regression, which has well known weaknesses, despite still being widely
17 used in ecology (Whittingham et al. 2006), random forest methods make no assumptions
18 about the distributions of variables and are robust to outliers in factors. They can also handle
19 situations where the number of factors exceeds the number of observations and have a novel
20 variable importance measure, which does not suffer the shortcomings of traditional variable
21 selection methods, such as selecting only one or two variables among a group of equally good
22 but highly correlated predictors (Cutler et al. 2007). Thus, the greater range of values of
23 explained variance in the RF outcomes compared to the STEP outcomes may be indicative of
24 their greater robustness and ability to distinguish meaningfulness relationships. Furthermore,
25 many studies that have applied classical regression approaches to understand the drivers of

1 forest cover changes (e.g. Jaimes et al. 2010; Gao and Li 2011; Freitas et al. 2013; Gong et al.
2 2013) may have had to use a restricted number of factors to be able to satisfy requirements of
3 normality, which could have hindered the analyses, whereas the flexibility and robustness of
4 RF overcomes such limitations.

5 Despite its advantages, RF used to be one main limitation. Unlike traditional
6 regression methods, RF did not produce relationships between independent and metrics that
7 have simple representations (such as linear equations), making ecological interpretation
8 difficult (Cutler et al. 2007). Nevertheless, the R "extendedForest" library has overcome this
9 issue. This package allows us to generate partial plots, which indicate the direction and form
10 of the independent response of a variable. Therefore, we can now convert the RF outcomes
11 into equations for quantitatively predicting changes in DF and FF metrics that might arise
12 from changes in the BGP and S-E factors considered here. Additionally, RF has exploited
13 structure in our high-dimensional data set not "visible" to STEP in the GRD and PD selected
14 models to provide an apparently clearer picture of these metrics' relationships to the factors.

15

16 **5 Conclusion**

17 Understanding spatial relationships between patterns of DF/FF metrics and S-E/BGP factors
18 is important for land use management. The main contribution of this study is the testing of a
19 relatively new application of RF for detecting this kind of relationship, its application to a
20 very large dataset, and its comparison with a traditional multiple linear regression method.
21 We found that RF performs better than multiple regression at explaining metrics describing
22 forest patch patterns (PD and ENN) and broader landscape structures (GRD and CA). Given
23 the well-established advantages of decision-tree-based methods over those of classical
24 multiple regression (Breiman et al. 1984; Breiman 2001; Prasad et al. 2006; Cutler et al.
25 2007; Cutler et al. 2008; Pitcher et al. 2011; Ellis et al. 2012; Cutler 2013; Smith et al. 2013),

1 we suggest that the reasons for these differences are likely to be because the patch-pattern
2 metrics and broader landscape structures vary in less smooth or monotonic ways (McGarigal
3 et al. 2012) – ways that RF is able to capture, but multiple regression is not. Still, we have
4 shown that RF provides a promising methodology for identifying these relationships, and that
5 it has the potential to be an effective tool for providing essential information for aiding land
6 use management decisions, not only in terms of planning, but also for conservation actions,
7 as proposed by Zanella et al. (2012), in cases of high rates of anthropogenic biodiversity loss,
8 as it is the case of the Atlantic Forest.

9 The initial investigation reported in the present study is, however, only a first step in
10 exploiting this method's potential. One aspect that requires further consideration is the scale
11 of the study area and the very wide variety of S-E and BGP contexts, which it encompasses.
12 Even in relatively small areas, a multitude of diverse factors are at work (Qasim et al. 2013),
13 and variations in contexts may have influenced model performance in the present study.
14 Landscape pattern is scale-sensitive (Gao and Li 2011) and the unusually large degree of
15 heterogeneity in the Atlantic forest biome is likely only to exacerbate this issue. Policies need
16 to be crafted at appropriate spatial scales and with specific contexts in mind. Thus, an
17 important development of this initial study of RF application to cases of DF and FF would be
18 to repeat it at different spatial scales, to identify more precisely the S-E and BGP factors
19 associated with these processes.

20

References

- Aide TM, Grau HR (2004) Globalization, Migration, and Latin American Ecosystems. *Science* (80-) 305:1915–1916. doi: 10.1126/science.1103179
- Apan AA, Peterson JA (1998) Probing tropical deforestation: the use of GIS and statistical
 25 analysis of georeferenced data. *Appl Geogr* 18:137–152.
- Barbier EB, Burgess JC, Grainger A (2010) The forest transition: Towards a more
 comprehensive theoretical framework. *Land use policy* 27:98–107. doi:
 10.1016/j.landusepol.2009.02.001
- Barton K (2014) Multi-model inference. *R Packag MuMIn* version 1105 46.
- 30 Beilin R, Lindborg R, Stenseke M, et al (2014) Analysing how drivers of agricultural land
 abandonment affect biodiversity and cultural landscapes using case studies from
 Scandinavia, Iberia and Oceania. *Land use policy* 36:60–72. doi:
 10.1016/j.landusepol.2013.07.003
- Bonilla-Moheno M, Aide TM, Clark ML (2012) The influence of socioeconomic,
 35 environmental, and demographic factors on municipality-scale land-cover change in
 Mexico. *Reg Environ Chang* 12:543–557. doi: 10.1007/s10113-011-0268-z
- Breiman L (2001) Random forests. *Mach Learn* 45:5–32.
- Breiman L, Friedman J, Stone CC, et al (1984) *Classification and regression trees*.
 Wadsworth, Belmont Calif.
- 40 Burnham K, Anderson D (2002) *Model selection and multimodel inference: a practical
 information-theoretic approach*, 2nd ed. New York
- Butler BJ, Swenson JJ, Alig RJ (2004) Forest fragmentation in the Pacific Northwest:
 Quantification and correlations. *For Ecol Manage* 189:363–373. doi:
 10.1016/j.foreco.2003.09.013
- 45 Caldas MM, Goodin D, Sherwood S, et al (2013) Land-cover change in the Paraguayan

- Chaco: 2000–2011. *J Land Use Sci* 4248:1–18. doi: 10.1080/1747423X.2013.807314
- Carvalho LMT de, Scolforo JR (2008) Inventário Florestal de Minas Gerais: Monitoramento da Flora Nativa 2005-2007. Editora da UFLA, Lavras
- Cutler A (2013) *Trees and Random Forests*. NIH 1R15AG037392-01.
- 50 Cutler A, Cutler DR, Stevens JR (2008) Tree-based methods. In: Xiaochun Li, Xu R (eds) *High-Dimensional Data Analysis in Cancer Research*. Springer New York, pp 89–109
- Cutler A, Stevens J (2006) Random Forests for microarrays. *Methods Enzymol* 422–432.
- Cutler DR, Edwards Jr. TC, Beard KH, et al (2007) Random Forests for Classification in Ecology. *Ecology* 88:2783–2792. doi: 10.1890/07-0539.1
- 55 Dean W (1996) *With broadax and firebrand: the destruction of the Brazilian Atlantic Forest*. California
- DeFries RS, Foley JA, Asner GP (2004) Land-use choices: Balancing human needs and ecosystem function. *Front Ecol Environ* 2:249–257. doi: 10.1890/1540-9295(2004)002[0249:LCBHNA]2.0.CO;2
- 60 Efroymson M (1960) Multiple regression analysis. *Math methods Digit Comput* 191–203.
- Ellis N, Smith SJ, Pitcher CR (2012) Gradient forests: calculating importance gradients on physical predictors. *Ecology* 93:156–68.
- Farinaci JS, Batistella M (2012) Variação na cobertura vegetal nativa em São Paulo: um panorama do conhecimento atual. *Rev Árvore* 36:695–705. doi: 10.1590/S0100-67622012000400011
- 65 Fearnside P (2016) *The Roles and Movements of Actors in the Deforestation of Brazilian Amazonia*. doi: Artn 23
- Ferreira MP, Alves DS, Shimabukuro YE (2015) Forest dynamics and land-use transitions in the Brazilian Atlantic Forest: the case of sugarcane expansion. *Reg Environ Chang*
- 70

15:365–377. doi: 10.1007/s10113-014-0652-6

Flamenco-Sandoval A, Martínez Ramos M, Masera OR (2007) Assessing implications of land-use and land-cover change dynamics for conservation of a highly diverse tropical rain forest. *Biol Conserv* 138:131–145. doi: 10.1016/j.biocon.2007.04.022

75 Foley JA (2005) Global Consequences of Land Use. *Science* (80-) 309:570–574. doi: 10.1126/science.1111772

Forman RTT, Alexander LE (1998) Roads and Their Major Ecological Effects. *Annu Rev Ecol Syst* 29:207–231. doi: 10.1146/annurev.ecolsys.29.1.207

Freitas MWD, Santos JR Dos, Alves DS (2013) Land-use and land-cover change processes in
80 the Upper Uruguay Basin: linking environmental and socioeconomic variables. *Landsc Ecol* 28:311–327. doi: 10.1007/s10980-012-9838-9

Freitas SR, Hawbaker TJ, Metzger JP (2010) Effects of roads, topography, and land use on forest cover dynamics in the Brazilian Atlantic Forest. *For Ecol Manage* 259:410–417. doi: 10.1016/j.foreco.2009.10.036

85 Fu W, Liu S, Degloria SD, et al (2010) Characterizing the “fragmentation–barrier” effect of road networks on landscape connectivity: A case study in Xishuangbanna, Southwest China. *Landsc Urban Plan* 95:122–129. doi: 10.1016/j.landurbplan.2009.12.009

Furundzic D (1998) Application example of neural networks for time series analysis : Rainfall — runoff modeling. 64:383–396.

90 Gao J, Li S (2011) Detecting spatially non-stationary and scale-dependent relationships between urban landscape fragmentation and related factors using Geographically Weighted Regression. *Appl Geogr* 31:292–302. doi: 10.1016/j.apgeog.2010.06.003

Gaughan AE, Binford MW, Southworth J (2009) Tourism, forest conversion, and land transformations in the Angkor basin, Cambodia. *Appl Geogr* 29:212–223. doi:
95 10.1016/j.apgeog.2008.09.007

Geist HJ, Lambin EF (2001) What drives tropical deforestation? LUCR Report Series No. 4.

Louvain-la-Neuve

Geist HJ, Lambin EF (2002) Proximate Causes and Underlying Driving Forces of Tropical

Deforestation. *Bioscience* 52:143–150. doi: 10.1641/0006-

100 3568(2002)052[0143:PCAUDF]2.0.CO;2

Gilbert A, Chakraborty J (2011) Using geographically weighted regression for environmental

justice analysis: Cumulative cancer risks from air toxics in Florida. *Soc Sci Res* 40:273–

286. doi: 10.1016/j.ssresearch.2010.08.006

Gislason PO, Benediktsson JA, Sveinsson JR (2006) Random Forests for land cover

105 classification. *Pattern Recognit Lett* 27:294–300. doi: 10.1016/j.patrec.2005.08.011

Gong C, Yu S, Joesting H, Chen J (2013) Determining socioeconomic drivers of urban forest

fragmentation with historical remote sensing images. *Landsc Urban Plan* 117:57–65.

doi: 10.1016/j.landurbplan.2013.04.009

Hecht SB, Kandel S, Gomes I, et al (2006) Globalization, forest resurgence, and

110 environmental politics in El Salvador. *World Dev* 34:308–323. doi:

10.1016/j.worlddev.2005.09.005

Hocking RR, Mar N (1976) A Biometrics Invited Paper . The Analysis and Selection of

Variables in Linear Regression. 32:1–49.

IBGE (2017) Estados. Minas Gerais. <http://www.ibge.gov.br/estadosat/perfil.php?sigla=mg>.

115 Accessed 15 Jun 2017

Jaimes NBP, Bosque Sendra J, Franco R, et al (2010) Exploring the driving forces behind

deforestation in the state of Mexico (Mexico) using geographically weighted regression.

Appl Geogr 30:576–591. doi: 10.1016/j.apgeog.2010.05.004

James G, Witten D, Hastie T, Tibshirani R (2013) An Introduction to Statistical Learning:

120 with Applications in R. New York

- Kadane JB, Lazar N a (2004) Methods and Criteria for Model Selection. *J Am Stat Assoc* 99:279–290. doi: 10.1198/016214504000000269
- Kauano ÉE, Torezan JMD, Cardoso FCG, Marques MCM (2012) Landscape structure in the northern coast of Paraná state, a hotspot for the brazilian Atlantic Forest conservation. *Rev Árvore* 36:961–970. doi: 10.1590/S0100-67622012000500018
- 125
- Killeen TJ, Calderon V, Soria L, et al (2007) Thirty years of land-cover change in Bolivia. *Ambio* 36:600–6. doi: 10.1579/0044-7447(2007)36[600:TYOLCI]2.0.CO;2
- Kleinbaum D, Kupper L, Nizam A, Rosenberg E (1998) Applied regression analysis and other multivariable methods, third. Duxbury Press., Pacific Grove - CA
- 130
- Lambin EF, Geist HJ (2006) Land-Use and Land-Cover Change: Local Processes and Global Impacts, 1st edn. Springer Verlag, Berlin
- Laurance WF, Cochrane MA, Bergen S, et al (2001) The Future of the Brazilian Amazon. *Science* (80-) 291:438–439. doi: 10.1126/science.291.5503.438
- Leal CGCG, Pompeu PS, Gardner TA, et al (2016) Multi-scale assessment of human-induced changes to Amazonian instream habitats. *Landsc Ecol* 31:1725–1745. doi: 10.1007/s10980-016-0358-x
- 135
- Leinbach TR (1995) Transport and third world development: review, issues, and prescription. *Transp Res Part A Policy Pract* 29:337–344. doi: 10.1016/0965-8564(94)00035-9
- Lira PK, Ewers RM, Banks-Leite C, et al (2012) Evaluating the legacy of landscape history: extinction debt and species credit in bird and small mammal assemblages in the Brazilian Atlantic Forest. *J Appl Ecol* 49:1325–1333. doi: 10.1111/j.1365-2664.2012.02214.x
- 140
- Malhi Y, Gardner T a., Goldsmith GR, et al (2014) Tropical Forests in the Anthropocene. *Annu Rev Environ Resour* 39:125–159. doi: 10.1146/annurev-environ-030713-155141
- 145
- Margules CR, Pressey RL (2000) Systematic conservation planning. *Nature* 405:243–53. doi:

10.1038/35012251

McGarigal K, Cushman S, Ene E (2012) FragStats v4: Spatial Pattern Analysis Program for Categorical and Continuous Maps Computer software program produced by the authors at the University of Massachusetts.

150 <http://www.umass.edu/landeco/research/fragstats/fragstats.html>. Accessed 1 May 2016

Metzger JP, Casatti L (2006) Do diagnóstico à conservação da biodiversidade : o estado da arte do programa BIOTA / FAPESP. 6:1–23.

Oliveira-Filho A, Fontes M (2000) Patterns of Floristic Differentiation among Atlantic Forests in Southeastern Brazil and the Influence of Climate. *Biotropica* 32:793–810. doi:

155 10.1111/j.1744-7429.2000.tb00619.x

Oliveira VHF, Barlow J, Gardner T, Louzada J (2017) Do we select the best metrics for assessing land use effects on biodiversity? *Basic Appl Ecol*. doi:

10.1016/j.baae.2017.03.002

Parcerisas L, Marull J, Pino J, et al (2012) Land use changes, landscape ecology and their socioeconomic driving forces in the Spanish Mediterranean coast (El Maresme County, 1850–2005). *Environ Sci Policy* 23:120–132. doi: 10.1016/j.envsci.2012.08.002

160

Parés-Ramos IK, Gould WA, Aide TM (2008) Agricultural abandonment, suburban growth, and forest expansion in Puerto Rico between 1991 and 2000. *Ecol Soc*. doi: 1

Perz SG (2004) Are agricultural production and forest conservation compatible? Agricultural diversity, agricultural incomes and primary forest cover among small farm colonists in the Amazon. *World Dev* 32:957–977. doi: 10.1016/j.worlddev.2003.10.012

165

Perz SG, Caldas MM, Arima E, Walker RT (2007) Unofficial road building in the Amazon: Socioeconomic and biophysical explanations. *Dev Change* 38:529–551. doi:

10.1111/j.1467-7660.2007.00422.x

170 Pfaff ASP (1999) What Drives Deforestation in the Brazilian Amazon? Evidence from

- Satellite and Socioeconomic Data*. *J Environ Econ Manage* 37:26–43. doi:
10.1006/jeem.1998.1056
- Pitcher CR, Ellis N, Smith SJ (2011) Example analysis of biodiversity survey data with R package gradientForest Gradient Forest basics. 1–16.
- 175 Prasad AM, Iverson LR, Liaw A, et al (2006) Newer Tree Classification and Techniques :
Forests Random Prediction Bagging for Ecological Regression. *Ecosystems* 9:181–199.
doi: 10.1007/S10021-005-0054-1
- Qasim M, Hubacek K, Termansen M (2013) Underlying and proximate driving causes of
land use change in district Swat, Pakistan. *Land use policy* 34:146–157. doi:
180 10.1016/j.landusepol.2013.02.008
- Quezada ML, Arroyo-Rodríguez V, Pérez-Silva E, Aide TM (2013) Land cover changes in
the Lachuá region, Guatemala: patterns, proximate causes, and underlying driving forces
over the last 50 years. *Reg Environ Chang* 14:1139–1149. doi: 10.1007/s10113-013-
0548-x
- 185 Redo DJ, Aide TM, Clark ML (2012) The Relative Importance of Socioeconomic and
Environmental Variables in Explaining Land Change in Bolivia, 2001–2010. *Ann Assoc
Am Geogr* 102:778–807. doi: 10.1080/00045608.2012.678036
- Ribeiro MCM, Metzger JPJP, Martensen AC, et al (2009) The Brazilian Atlantic Forest:
How much is left, and how is the remaining forest distributed? Implications for
190 conservation. *Biol Conserv* 142:1141–1153. doi: 10.1016/j.biocon.2009.02.021
- Richards PD, Walkerb RT, Arima EY (2008) NIH Public Access. *Glob Env Chang* 144:724–
732. doi: 10.1038/jid.2014.371
- Riitters KH, Neil RVO, Hunsaker CT, et al (1995) A factor analysis of landscape pattern and
structure metrics. *Landsc Ecol* 10:23–39.
- 195 Scolforo J, Carvalho LMT de (2006) Mapeamento e inventário da flora nativa e dos

reflorestamentos de Minas Gerais, 2nd edn. UFLA

Scolforo JR, Oliveira AD de, Carvalho LMT de (2008) Zoneamento ecológico-econômico do estado de minas gerais: Componente sócioeconômico. UFLA, IAVRAS

200 Silva WG, Metzger JP, Simões S, Simonetti C (2007) Relief influence on the spatial distribution of the Atlantic Forest cover on the Ibiúna Plateau, SP. *Brazilian J Biol* 67:403–11.

Smith PF, Ganesh S, Liu P (2013) A comparison of random forest regression and multiple linear regression for prediction in neuroscience. *J Neurosci Methods* 220:85–91. doi: 10.1016/j.jneumeth.2013.08.024

205 Smith SJ, Ellis N, Pitcher CR (2011) Conditional variable importance in R package `extendedForest`.

SOS Mata Atlântica/INPE (2014) Atlas dos remanescentes de Mata Atlântica período 2012–2013. São Paulo, Brazil

210 Teixeira AMG, Soares-Filho BS, Freitas SR, Metzger JP (2009) Modeling landscape dynamics in an Atlantic Rainforest region: Implications for conservation. *For Ecol Manage* 257:1219–1230. doi: 10.1016/j.foreco.2008.10.011

Venables WN, Ripley BD (2002) *Modern Applied Statistics with S*, Fourth edi. Springer-Verlag, New York

215 Wei C-L, Rowe GT, Escobar-Briones E, et al (2010) Global patterns and predictions of seafloor biomass using random forests. *PLoS One* 5:e15323. doi: 10.1371/journal.pone.0015323

Whittingham MJ, Stephens P a, Bradbury RB, Freckleton RP (2006) Why do we still use stepwise modelling in ecology and behaviour? *J Anim Ecol* 75:1182–9. doi: 10.1111/j.1365-2656.2006.01141.x

220 Wright SJ, Samaniego MJ (2008) Historical, demographic, and economic correlates of land-

use change in the Republic of Panama.

Zanella L, Borém R, Souza C, et al (2012) Atlantic Forest Fragmentation Analysis and Landscape Restoration Management Scenarios. *Nat ...* 10:57–63.

225 **Table 1** Descriptions of deforestation (DF) and forest fragmentation (FF) metrics (dependent variables)

Metric	Category	Formulae	Description (unit) ^a
Growth rate of deforestation (GRD)	Rate of deforestation	$GRD = \frac{(Df - Di) / Di}{t}$	<p>Growth rate of deforestation from 2003 to 2011. Di = Total area deforested in 2003. Df = Total area deforested in 2011.</p> <p>t = number of years considered (in our case, eight years).</p> <p>Percentage.</p>
Mean Euclidean Nearest-Neighbour (ENN)	Forest patch isolation	$ENN = \frac{\sum_{j=1}^n h_{ij}}{n_i}$	<p>ENN equals the mean distance to the nearest neighbouring patch of forest, based on shortest edge-to-edge distance. h_{ij} = distance (m) from patch j to nearest neighbouring patch</p>

			of the same type (<i>i</i> , in this case forest). n_i = number of patches of cover type <i>i</i> (forest).
Total remaining forest (CA)	Remaining forest quantification	$CA = A \left(\frac{1}{10,000} \right)$	CA equals the total area (m ²) of the landscape, divided by 10,000 (to convert to hectares). <i>A</i> = total landscape area (m ²). CA is important because it defines the extent of the landscape.
Patch density (PD)	Forest spatial structure	$PD = \frac{n_i}{A} (10,00000)$	Patch density increases with a greater number of patches within a reference area and therefore reflects landscape fragmentation.

^aDetails can be found in Mcgarigal et al. (2012).

Fig 1 Minas Gerais State, BR and the 518 municipalities used in this study. The inset map on
230 the left show the location of Minas Gerais State within Brazil

Fig 2 Relative importance plot for factors from random forest (RF) and stepwise multiple
regression (STEP) analysis approaches, in percentage (%). Factors are defined in Table S2
(supplementary material ESM2). The eight selected models from both approaches are: the
235 growth rate of deforestation – RF-GRD and STEP-GRD; the total remaining forest – RF-CA
and STEP-CA; the mean Euclidean nearest-neighbour distance – RF-ENN and STEP-ENN;
and the patch density – RF-PD and STEP-P. Metrics used in models are defined in Table 1.
Note that each model shows only the most important predictors. * Percentage of variance
explained in each model

240

Fig 3 Partial contribution of socio-economic (S-E) and bio-geophysical (BGP) factors to
deforestation (DF) and forest fragmentation (FF) in Minas Gerais, Brazil, derived from RF
analysis approach. Factors are defined in Table S2 (Supplementary material ESM2). A) The
growth rate of deforestation (RF-GRD); B) Patch density (RF-PD); C) The total remaining
245 forest area (RF-CA); and D) The Euclidean nearest-neighbour distance (RF-ENN). Metrics
used in models are defined in Table 1