

SITS: Data Analysis and Machine Learning using Satellite Image Time Series

Rolf Simoes *National Institute for Space Research (INPE), Brazil*

Gilberto Camara *National Institute for Space Research (INPE), Brazil*

Alexandre Carvalho *Institute for Applied Economics Research (IPEA), Brazil*

Victor Maus *International Institute for Applied System Analysis (IIASA), Switzerland*

Gilberto Queiroz *National Institute for Space Research (INPE), Brazil*

Using time series derived from big Earth Observation data sets is one of the leading research trends in Land Use Science and Remote Sensing. One of the more promising uses of satellite time series is its application for classification of land use and land cover, since our growing demand for natural resources has caused major environmental impacts. Here, we present the open source R package for satellite image time series analysis, the `sits` package. The `sits` provides support on how to use statistical learning techniques with image time series. These methods include linear and quadratic discrimination analysis, support vector machines, random forests and neural networks.

Introduction

Earth observation satellites provide a continuous and consistent set of information about the Earth's land and oceans. Most space agencies have adopted an open data policy, making unprecedented amounts of satellite data available for research and operational use. This data deluge has brought about a major challenge: *How to design and build technologies that allow the Earth observation community to analyse big data sets?*

The approach taken in the current work is to develop data analysis methods that work with satellite image time series. The time series are obtained by taking calibrated and comparable measures of the same location in Earth at different times. These measures can be obtained by a single sensor (*e.g.*, MODIS) or by combining different sensors (*e.g.*, LANDSAT-8 and SENTINEL-2). If obtained by frequent revisits, the temporal resolution of these data sets can capture the most important land use changes.

Time series of remote sensing data show that land cover changes do not always occur in a progressive and gradual way, but they may also show periods of rapid and abrupt change followed either by a quick recovery [Lambin et al., 2003]. Analyses of multiyear time series of land surface attributes, their fine-scale spatial pattern, and their seasonal evolution leads to a broader view of land-cover change. Satellite image time series have already been applied to applications such as mapping for detecting forest disturbance [Kennedy et al., 2010], ecology dynamics [Pasquarella et al., 2016], agricultural intensification [Galford et al., 2008] and its impacts on deforestation [Arvor et al., 2012].

In this paper, we present an open source R package for satellite image time series analysis `sits`. The `sits` package provides support on how to use statistical learning techniques with image time series. In a broad sense, statistical learning refers to a class of algorithms for classification and regression analysis [Hastie et al., 2009]. These methods include linear and quadratic discrimination analysis, support vector machines, random forests and neural networks. In a typical classification problem, we have measures that capture class attributes. Based on these measures, referred as training data, one's task is to select a predictive model that allows inferring classes of a larger data set.

In what follows, we describe the main characteristics of the `sits`. The first part describes the basic data structures used in it and the tools used for visualisation and data exploration. Then we show how to do data acquisition from external sources, with an emphasis on the WTSS (an acronym for Web Time Series Service) [Queiroz et al., 2015]. The next sections describe filtering and clustering techniques. We then discuss machine learning techniques for satellite image time series data and how to apply them to image time series. Finally, we present validation methods.

Data Handling and Visualisation Basics in `sits`

The basic data unit in the `sits` package is the “`sits tibble`”, which is a way of organizing a set of time series data with associated spatial information. In R, a `tibble` differs from the traditional data frame, insofar as a `tibble` can contain lists embedded as column arguments. Tibbles are part of the `tidyverse`, a collection of R package designed to work together in data manipulation. The `tidyverse` includes packages such as `ggplot2`, `dplyr` and `purrr` [Wickham and Grolemund, 2017]. The `sits` makes extensive use of the `tidyverse`.

For a better explanation of how the “`sits tibble`” works, we will read a data set containing 2,115 labelled samples of land cover in Mato Grosso state of Brazil. This state has 903,357 km² of extension, being the third largest state of Brazil. It includes three of Brazil's biomes: Amazonia, Cerrado and Pantanal. It is the most important agricultural frontier of Brazil and is Brazil's largest producer of soybeans, corn and cotton.

The samples contain time series extracted from the MODIS MOD13Q1 product from NASA from 2000 to 2016, provided every 16 days at 250-meter spatial resolution in the Sinusoidal projection. Based on ground surveys and high resolution imagery, we selected 2,115 samples of nine classes: forest, cerrado, pasture, soybean-fallow, fallow-cotton, soybean-cotton, soybean-corn, soybean-millet, and soybean-sunflower.

```
# data set of samples
# print the first three samples
samples_MT_9classes[1:3,]
```

```
## # A tibble: 3 x 7
##   longitude latitude start_date end_date  label  coverage  time_series
##   <dbl>    <dbl> <date>    <date>    <chr>    <chr>    <list>
## 1   -55.2   -10.8 2013-09-14 2014-08-29 Pasture mod13q1_512 <tibble [2~
## 2   -57.8    -9.76 2006-09-14 2007-08-29 Pasture mod13q1_512 <tibble [2~
```

```
## 3      -51.9   -13.4  2014-09-14 2015-08-29 Pasture mod13q1_512 <tibble [2~
```

The “sits tibble” contains data and metadata. The first six columns contain the metadata: spatial and temporal location, label assigned to the sample, and coverage from where the data has been extracted. The spatial location is given in longitude and latitude coordinates for the “WGS84” ellipsoid. For example, the first sample has been labelled “Pasture”, at location (-55.1852, -10.8387), and is considered valid for the period (2013-09-14, 2014-08-29). Informing the dates where the label is valid is crucial for correct classification. In this case, the researchers involved in labelling the samples chose to use the agricultural calendar in Brazil, where the spring crop is planted in the months of September and October, and the autumn crop is planted in the months of February and March. For other applications and other countries, the relevant dates will most likely be different from those used in the example.

The “sits tibble” also contains the time series data for each spatiotemporal location. The timeseries data is also organized as a tibble, with a column with the dates and the other columns with the values for each spectral band.

```
# print the first 10 time series records of the first sample
samples_MT_9classes$time_series[[1]][1:3,]
```

```
## # A tibble: 3 x 7
##   Index      ndvi   evi   nir   mir   blue   red
##   <date>    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 2013-09-14 0.424 0.280 0.288 0.244 0.0605 0.116
## 2 2013-09-30 0.467 0.264 0.257 0.167 0.0357 0.0933
## 3 2013-10-16 0.504 0.299 0.266 0.202 0.0405 0.0877
```

The sits package provides functions for data manipulation and displaying information of a “sits tibble”. For example, the command `sits_labels()` that shows the labels of the sample set and their frequencies.

```
sits_labels(samples_MT_9classes)
```

```
## # A tibble: 9 x 3
##   label      count  freq
##   <chr>    <int> <dbl>
## 1 Cerrado      400 0.189
## 2 Fallow_Cotton  34 0.0161
## 3 Forest      138 0.0652
## 4 Pasture     370 0.175
## 5 Soy_Corn     398 0.188
## 6 Soy_Cotton   399 0.189
## 7 Soy_Fallow    88 0.0416
## 8 Soy_Millet   235 0.111
## 9 Soy_Sunflower  53 0.0251
```

In many cases, it is useful to relabel the data set. For example, there may be situations when one wants to use a smaller set of labels, since samples in one label on the original set may not be distinguishable from samples with other labels. We then should use `sits_relabel()`, which requires a conversion list (for details, see `?sits_relabel`).

Given that we have used the tibble data format for the metadata and the embedded time series, one can use the functions of the `dplyr`, `tidyr` and `purrr` packages of the tidyverse [Wickham and Grolemund, 2017] to process the data. For example, the following code uses the `sits_select()` function to get a subset of the sample data set with two bands ("ndvi" and "evi") and then uses the `dplyr::filter()` function to select the samples labelled either as "Cerrado" or "Pasture". We can then use the `sits_plot()` function to display the time series. Given a small number of samples to display, the `sits_plot()` function tries to group as many spatial locations together. In the following example, the first 15 samples of the "Cerrado" class all refer to the same spatial location in consecutive time periods. For this reason, these samples are plotted together.

```
# select the "ndvi" bands
samples_ndvi.tb <-
  sits_select(samples_MT_9classes,
              bands = c("ndvi"))
# select only the samples with the cerrado label
samples_cerrado.tb <-
  dplyr::filter(samples_ndvi.tb,
                label == "Cerrado")
# plot the first 15 samples (different dates for the same points)
sits_plot(samples_cerrado.tb[1:15,])
```

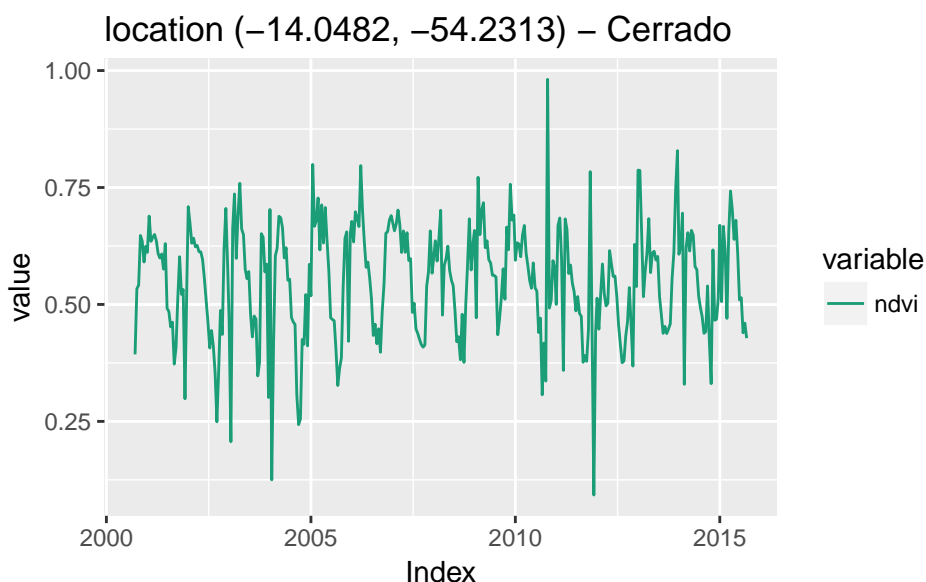


Figure 1: Plot of the first 15 'Cerrado' samples from data set `samples_MT_9classes` (different dates for the same point location).

For a large number of samples, where the amount of individual plots would be substantial, the default visualisation combines all samples together in a single temporal interval (even if they are valid for different years). Therefore, all samples of the same band and the same label are aligned to a common interval. This plot is useful to show the spread of values for the time series of each band. The strong red line in the plot shows the median of the values, and the two orange lines are the first and third interquartile ranges. The `sits_plot()` function has different ways of working. Please, refer to the package documentation for more details.

```
# plot all cerrado samples together (shows the distribution)
sits_plot(samples_cerrado.tb)
```

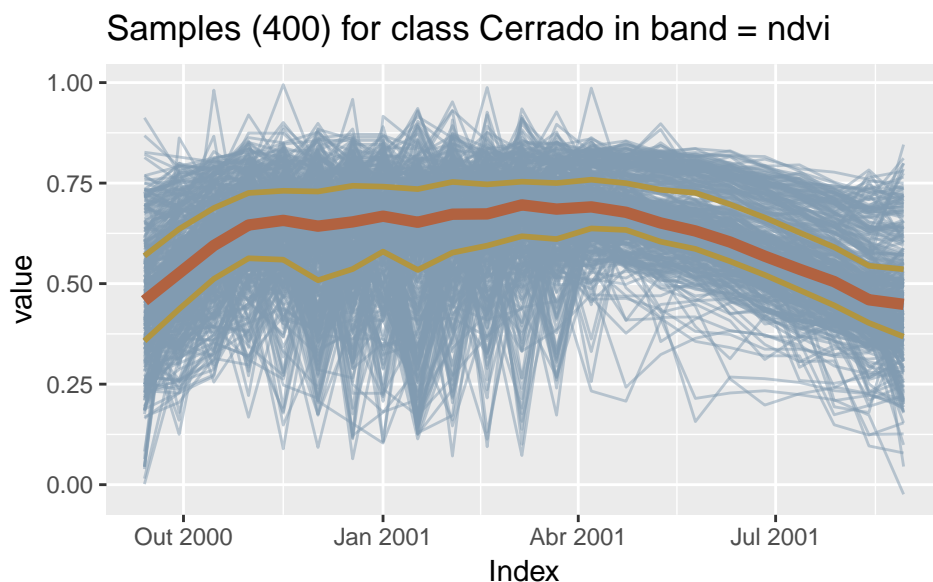


Figure 2: Plot of all 'Cerrado' samples from data set `samples_MT_9classes`.

Usually, samples are provided by experts whose take *in-loco* field observations or recognises land classes through high resolution images. In any case, we need access to a data source to fetch time series data regarding a spatiotemporal location of interest. The process of importing data samples is discussed in the next section.

Importing Data into `sits`

The `sits` package allows different methods of data input, including: (a) obtain data from a WTSS (Web Series Time Service); (b) obtain data from the SATVEG service developed by EMBRAPA (Brazil's Agriculture Research Agency). (c) read data stored in a time series in the ZOO format [Zeileis and Grothendieck, 2005]; (d) read a time series from a RasterBrick [Hijmans, 2015]. Option (d) will be described in the section where we describe raster processing. The WTSS service is a light-weight service, designed to retrieve time series for selected locations and periods [Vinhas et al., 2016], been implemented by the research team of the National Institute for Space Research to allow remote access to

time series data. The SATVEG service provides NDVI and EVI time series vegetation indices from MODIS image from whole Brazilian territory [EMBRAPA, 2014]. To view service details, the user needs to call `sits_services()` that provides information on the coverages available on the server.

After finding out which coverages are available at the different time series services, one may request specific information on each coverage by using `sits_coverage()`. This lists the contents of the data set, including source, bands, spatial extent and resolution, time range, and temporal resolution. This information is then stored in a tibble for later use.

```
# get information about a specific coverage from WTSS
coverage.tb <-
  sits_coverage(service = "WTSS",
                product  = "MOD13Q1",
                coverage = "mod13q1_512")
coverage.tb[, c("xmin", "xmax", "ymin", "ymax",
               "start_date", "end_date")]
```

```
## # A tibble: 1 x 6
##   xmin xmax ymin ymax start_date end_date
##   <dbl> <dbl> <dbl> <dbl> <date>    <date>
## 1 -81.2 -31.9 -30.0  10.0 2000-02-18 2017-02-18
```

The user can request one or more time series points using `sits_getdata()`. This function provides a general means of access to image time series. In its simplest fashion, the user provides the latitude and longitude of the desired location, the product and coverage names, the bands, and the start date and end date of the time series. If the start and end dates are not provided, all available period is retrieved. The result is a tibble that can be visualised using `sits_plot()`.

```
# a point in the transition forest pasture in Northern MT
# obtain a time series from the WTSS server for this point
series.tb <-
  sits_getdata(longitude = -55.57320,
               latitude  = -11.50566,
               service   = "WTSS",
               product    = "MOD13Q1",
               coverage   = "mod13q1_512",
               bands      = c("ndvi", "evi"),
               start_date = "2001-01-01",
               end_date   = "2016-12-31")
# plot the series
sits_plot(series.tb)
```

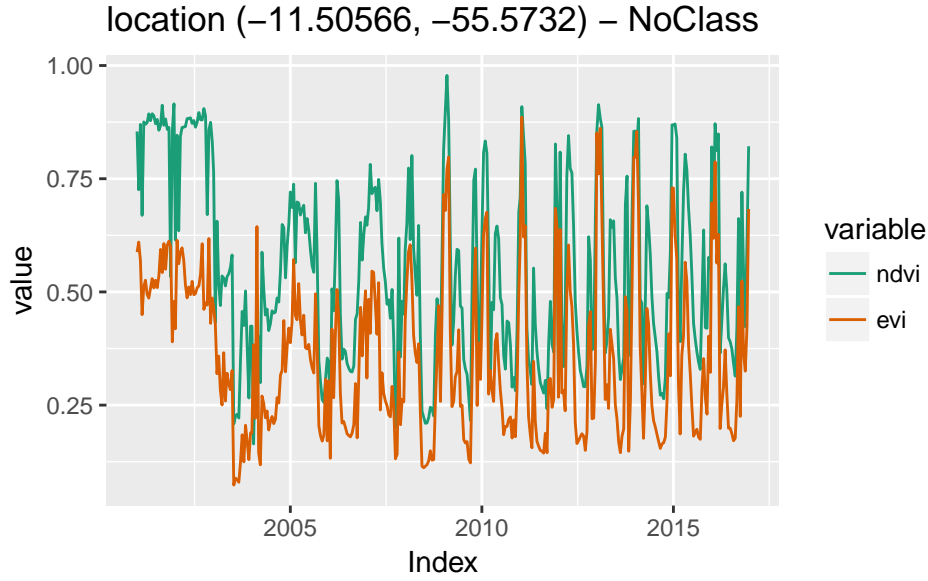


Figure 3: NDVI and EVI time series fetched from WTSS service.

A useful case is when users have a set of labelled samples, that are to be used as a training data set. In this case, one usually has trusted observations which are labelled and commonly stored in plain text CSV files. The `sits_getdata()` function can receive a CSV file path as an argument. The CSV file must provide for each time series, its latitude and longitude, the start and end dates, and a label associated to a ground sample.

After importing the samples time series, it is useful to explore the data and see how is it underlying structured and its inter-class separability. For example, We can note in the figure above the variability of 400 time series samples along time. Those samples were collected from different years and/or locations. The scattering behaviour is intrinsic to remote sensing data. Atmospheric noise, sun angle, interferences on observations or different equipments specifications, as well as the very nature of the climate-land dynamics can be sources of such variability [Atkinson et al., 2012]. One helpful technique to explore such properties is *cluster analysis*. In the following section we present a cluster technique supported by `sits`.

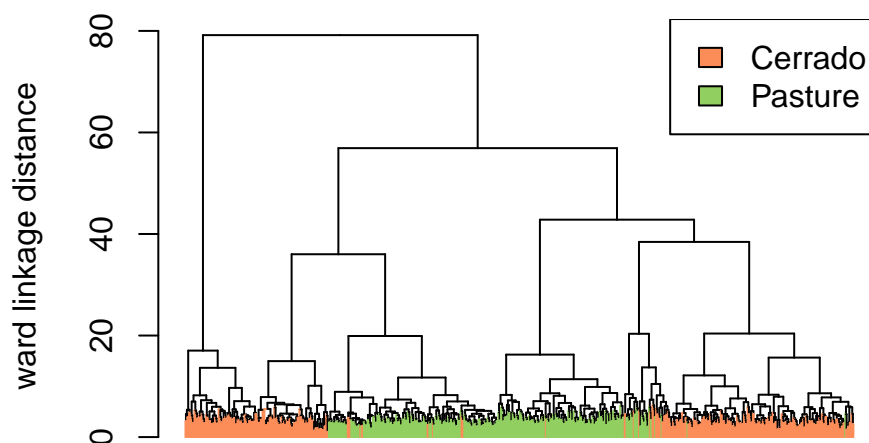
Clustering in satellite image time series

Cluster analysis has been used for many purposes in satellite image time series literature ranging from unsupervised classification [Petitjean et al., 2011], and pattern detection [Romani et al., 2011]. Here, we are interested in the second use of clustering, as a way to improve training data to use in machine learning classification models. In this regard, cluster analysis can assist the identification of structural *time series patterns*, and anomalous samples [Romani et al., 2011], [Chandola et al., 2009]. `sits` provides support for the agglomerative hierarchical clustering (AHC).

Hierarchical clustering is a family of methods that groups elements using a distance function to associate a real value to a pair of elements. From this distance measure, we can compute the dissimilarity between any two elements from the data set. Depending on the distance functions and linkage criteria, the algorithm decides which two clusters are merged at each iteration. The AHC approach is suitable for the purposes of samples data exploration awe its visualisation power and ease of use [Keogh et al., 2003]. Moreover, AHC does not require a predefined number of clusters as an initial parameter. This is an important feature in satellite image time series clustering since it is not easy to define the number of clusters present in a set of multi-attribute time series [Aghabozorgi et al., 2015].

The main result of AHC method is the *dendrogram*. A *dendrogram* is the ultrametric relation formed by the successive merges in the hierarchical process that can be represented by a tree. Dendrograms are quite useful to decide on the number of clusters has the data. It shows the height where each merging happened, which corresponds to the minimum distance between two clusters defined by a *linkage criterion*. The most common linkage criteria are: *single-linkage*, *complete-linkage*, *average-linkage*, and *Ward-linkage*. Complete-linkage prioritises the within-cluster dissimilarities, producing clusters with shorter distance samples. Complete-linkage clustering can be sensitive to outliers, that can increase the resulting intracluster data variance. As an alternative, Ward proposes a criteria to minimise the data variance by means of either *sum-of-squares* or *sum-of-squares-error* [Ward, 1963]. Ward's intuition is that clusters of multivariate observations, such as time series, should be approximately elliptical in shape [Hennig, 2015]. In sits, a dendrogram can be generated by `sits_dendrogram()`. The following codes illustrate how to create, visualise, and cut a dendrogram (for details, see `?sits_dendrogram()`).

```
# take a set of patterns for 2 classes
# create a dendrogram object with default clustering parameters
dendro <- sits_dendrogram(cerrado_2classes)
# plot the resulting dendrogram
sits_plot_dendrogram(cerrado_2classes,
                     dendro)
```



After the creation of a dendrogram, an important question emerges: *where to cut the dendrogram?* The answer depends on what are the purposes of the cluster analysis [Hennig, 2015]. If one is interested in an unsupervised classification, it is common to use *internal validity indices*, such as Silhouettes [Rousseeuw, 1987], to help determine the best number of clusters. However, if one is interested in understanding the structure of a labeled data set, or in the identification of sample anomaly, as we are here, one can recur to *external validity indices* to assist the semisupervised procedure that achieves the optimal correspondence between the clusters and classes partitions. In this regard, we need to balance two objectives: get clusters as large as possible, and get clusters as homogeneous as possible with respect to its known classes. To help this process, `sits` provides `sits_dendro_bestcut()` function that computes the external validity index *adjusted Rand index* (ARI) for a series of different number of generated clusters. The function returns the height where the cut of the dendrogram maximizes the index.

```
# search for the best height to cut the dendrogram
sits_dendro_bestcut(cerrado_2classes,
                    dendro)
```

```
##          k    height
## 6.000000 20.39655
```

This height optimises the ARI and generates 6 clusters. The ARI considers any pair of distinct samples and computes the following counts: a) the number of distinct pairs whose samples have the same label and are in the same cluster; b) the number of distinct pairs whose samples have the same label and are in different clusters; c) the number of distinct pairs whose samples have different labels and are in the same cluster; d) the number of distinct pairs whose samples have the different labels and are in different clusters. Here, *a* and *d* consists in all agreements, and *b* and *c* all disagreements. The ARI is obtained by

$$ARI = \frac{a + d - e}{a + d + b + c - e},$$

where *e* is the expected agreement, a random chance correction [Hubert and Arabie, 1985]. Different from others validity index such as Jaccard and Fowlkes-Mallows indices, ARI is more appropriate either when the number of clusters is outweighed by the number of labels (or conversely) and the number of samples is imbalanced between labels and clusters (which is usually the case).

```
# create 6 clusters by cutting the dendrogram at the linkage distance 20.39655
clusters.tb <-
  sits_cluster(cerrado_2classes,
               dendro,
               k = 6)
# show clusters samples frequency
sits_cluster_frequency(clusters.tb)
```

```
##
```

```
##           1    2    3    4    5    6 Total
## Cerrado 203  13  23  80    1  80   400
## Pasture   2 176  28    0 140    0   346
## Total   205 189  51  80 141  80   746
```

Note in this example that almost all clusters has a predomination of either “Cerrado” or “Pasture” classes with the exception of cluster 3. The contingency table plotted by `sits_cluster_frequency` is useful to identify anomalies in samples in a qualitative fashion. If one considers such cases as outliers, it is possible to exclude them using the function `sits_cluster_cleaner()` that removes all samples whose label (cluster) counts less than a given percentual relative to its cluster (label).

```
# clear sample outliers relative to clusters
# (those with less than 5% relative to intracluster)
clusters2.tb <- sits_cluster_cleaner(clusters.tb)
# show clusters samples frequency
sits_cluster_frequency(clusters2.tb)
```

```
##
##           1    2    3    4    5    6 Total
## Cerrado 203  13  23  80    0  80   399
## Pasture   0 176  28    0 140    0   344
## Total   203 189  51  80 140  80   743
```

The dendrogram can be used to evaluate cluster separability in a qualitative fashion, and helps the practitioner to determine where to cut the tree to obtain the final clusters [Liao, 2005].

Filtering techniques

Satellite image time series will always be contaminated by atmospheric influence, geolocation error, and directional effects [Lambin and Linderman, 2006]. In tropical regions, clouds are a major factor on reducing the quality of the data. Inter-annual climate variability also changes the phenological cycles of the vegetation, resulting in time series whose periods and intensities do not match on an year to year basis [Atkinson et al., 2012]. As we can note in The same land use and land cover class is subject to intrinsic year to year variability. To make the best use of available satellite data archives, methods for satellite image time series analysis need to deal with data sets that are *noisy* and *non-homogeneous*.

Several techniques have been used to correct and smooth such time-series vegetation index data, and to support the estimation of phenological parameters. These methods include curve fitting [Bradley et al., 2007], asymmetric Gaussian functions [Jonsson and Eklundh, 2002-08], wavelet decomposition [Sakamoto et al., 2005], the Savitzky–Golay filter [Chen et al., 2004], and the Whittaker smoother [Atzberger and Eilers, 2011].

For this purpose, we use data produced by the PRODES project [INPE, 2017] to extract our samples data set. The samples were extracted by the expert Dr. Rodrigo Bergotti (INPE). The PRODES project uses visual interpretation to detect new areas of clear cuts in the Amazon forest. Its interpreters are highly skilled and the accuracy of PRODES has been essential for Brazil's actions to control Amazon deforestation [Rajão and Hayes, 2009]. Our time series were generated from a combination of MODIS and Landsat 8 images covering the region corresponding to the standard Landsat Path/Row 226/064. This is an area in the East of the Pará state. The area has been chosen because the area is subject to strong cloud cover from November to March, which is a significant factor in degrading time series quality.

Each location information consists of its latitude and longitude, the start and end dates of an yearly time series, and the label, as showm below. Three land cover classes were identified: “Forest” (natural forest), “Deforestation_2014” and “Deforestation_2015” (areas identified by PRODES in years 2014 and 2015 as clear cuts, respectively) and “Pasture” (areas that had been deforested before 2014).

```
prodes_226_064[1:3,]
```

```
## # A tibble: 3 x 7
##   longitude latitude start_date end_date  label      coverage time_series
##   <dbl>    <dbl> <date>    <date>   <chr>      <chr>    <list>
## 1   -53.4    -6.53 2014-08-04 2015-07-22 Deforestation_2014 mixl8mod <tibble [23 x 4]>
## 2   -53.4    -6.51 2014-08-04 2015-07-22 Deforestation_2014 mixl8mod <tibble [23 x 4]>
## 3   -53.6    -6.11 2014-08-04 2015-07-22 Deforestation_2014 mixl8mod <tibble [23 x 4]>
```

The filtering algorithms provided by *sits* are: Savitzky–Golay filter (*sits_sgolay()*), Whitaker filter (*sits_whittaker()*), envelope filter (*sits_envelope()*) and, the “cloud filter” (*sits_cloud_filter()*). In what follows, we explain each technique and provide some examples using the afore described PRODES data set. To best visualize the effects of filtering we select only NDVI band of the time series.

Savitzky–Golay filter

The Savitzky–Golay filter works by fitting successive sub-sets of adjacent data points with a low-degree polynomial whose coefficients can be derived by through the method of linear least squares. The resulting i -th point of the filtered time series is obtained by substituting the polynomial independent terms by the n -adjacent points relative to the central point i in the original time series, as expressed in

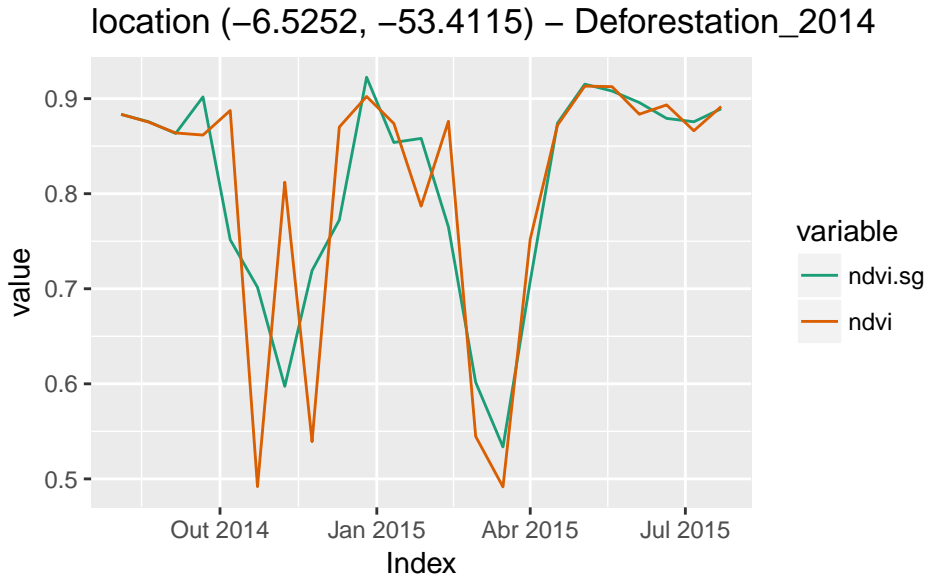
$$\hat{x}_i = \sum_{j=-n}^n C_j x_{i+j},$$

where \hat{x} is the the filtered time series, C_j is the polynomial coefficient respective to the j -th adjacent point and, x is the original time series. The coefficients have analytical solution when the time series is equally spaced and depends uniquely on polynomial

degree and the length of the sub set of adjacent data points (including the central point, $2n + 1$). More details on the procedure of coefficient determination can be found in [Savitzky and Golay \[1964\]](#).

The main purpose of Savitzky-Golay filter is to increase the signal-to-noise ratio with a minimum distortion of the time series. The difference between an original and a filtered time series using the Savitsky-Golay filter is shown in example bellow.

```
# Take the NDVI band of the first sample data set
point.tb <- sits_select(prodes_226_064[1,], bands = c("ndvi"))
# apply SavitzkyGolay filter
point_sg.tb <- sits_sgolay(point.tb)
# plot the series
sits_plot(sits_merge(point_sg.tb, point.tb))
```



Whittaker filter

The Whittaker smoother attempts to fit a curve that represents the raw data, but is penalized if subsequent points vary too much [\[Atzberger and Eilers, 2011\]](#). The Whittaker filter is a balancing between the residual to the original data and the “smoothness” of the fitted curve [\[de Rooi, 2013\]](#). The residual, as measured by the sum of squares of all n time series points deviations, in matrix form, is given by

$$RSS = (x - \hat{x})^T (x - \hat{x}),$$

where x and \hat{x} are the original and the filtered time series vectors, respectively. The “smoothness” is assumed to be the measure of the the sum of the squares of the third differences [\[Whittaker, 1922\]](#), which is given by

$$M_3 = (\hat{x}_4 - 3\hat{x}_3 + 3\hat{x}_2 - \hat{x}_1)^2 + (\hat{x}_5 - 3\hat{x}_4 + 3\hat{x}_3 - \hat{x}_2)^2 \\ + \dots + (\hat{x}_n - 3\hat{x}_{n-1} + 3\hat{x}_{n-2} - \hat{x}_{n-3})^2,$$

which in matrix form, can be expressed as

$$M_3 = \hat{x}^T D^T D \hat{x},$$

where

$$D = \begin{bmatrix} 1 & -3 & 3 & -1 & 0 & 0 & \cdots \\ 0 & 1 & -3 & 3 & -1 & 0 & \cdots \\ 0 & 0 & 1 & -3 & 3 & -1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

is the third order difference matrix.

Whittaker filter is obtained by finding a new time series \hat{x} whose points minimize the function

$$X(\hat{x}, \lambda) = RSS + \lambda M_3,$$

where λ , a scalar, works as an “smoothing wheight” parameter. The minimization can be obtained by differentiating the equivalent expression $(x - \hat{x})^T (x - \hat{x}) + \lambda \hat{x}^T D^T D \hat{x}$ with respect to \hat{x} and equating it to zero, whose solution gives the Whittaker filtered time series

$$\hat{x} = (\mathbf{I} + \lambda D^T D)^{-1} x.$$

The Whitakker filter can be a large but sparse optimisation problem, as we can note from D matrix. In `sits` implementation, the default λ parameter is 1.0. The example bellow depicts an original and a filtered time series.

```
# Take the NDVI band of the first sample data set
point.tb <- sits_select(prodes_226_064[1,], bands = c("ndvi"))
# apply Whittaker filter
point_whit.tb <- sits_whittaker(point.tb)
# plot the series
sits_plot(sits_merge(point_whit.tb, point.tb))
```

Envelope filter

This filter produces a bounding (superior or inferior) time series relative to an input signal. We can bound a time series signal by two basic operations: upper and lower dilations. Both operations assumes an unitary window to which the dilation occurs. An upper dilation can be defined as

$$u_i = \max_k (\{x_k : |k - i| \leq 1\}),$$

whereas an lower dilation is obtained by

$$l_i = \min_k (\{x_k : |k - i| \leq 1\}).$$

Here, x is the input time series and, k and i are vector indices.

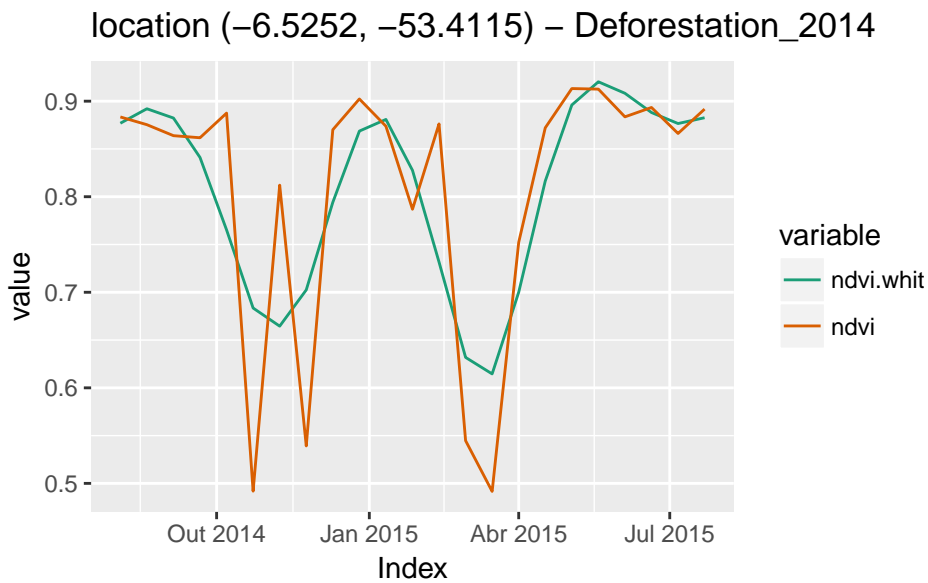
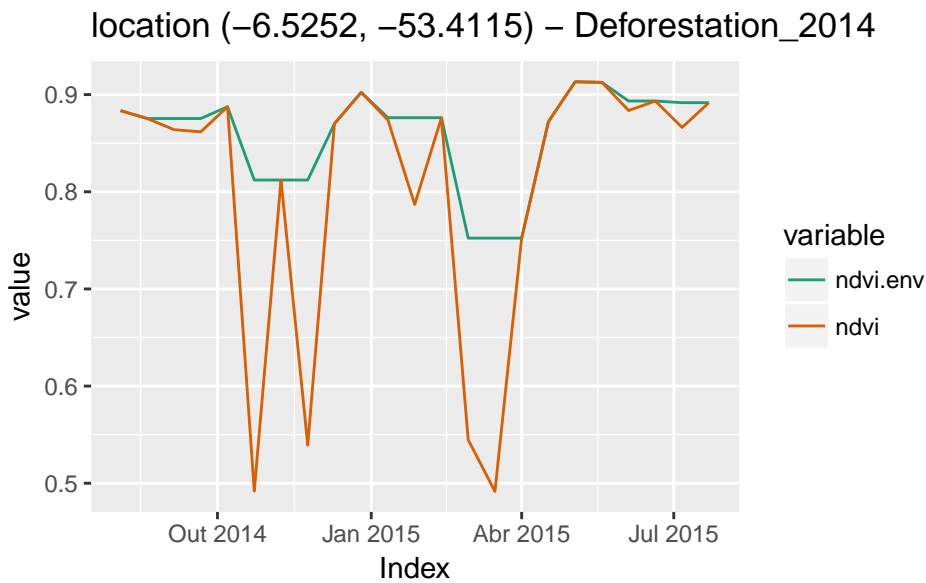


Figure 4: Whittaker filter applied on a one-year sample NDVI time series.

The envelope filter function can combine both upper and lower dilations recursively by a sequence of “U” (upper) and “L” (lower) characters passed as parameter. A repeated sequence of a same operation is equivalent to one operation with a larger dilation window. The final result of the recursive operations over the input time series is then returned.

This filter can be useful to remove isolated noises in the input signal. For example, an downward (upward) noise spike formed by one point in the series can be removed by an “UL” (“LU”) sequence of dilations. For spikes with two points, we can repeat each operation, for example “UULL” (the default sequence), to increase the dilation window and remove such imperfections on the time series. In the following example we can see an application of `sits_envelope()` function.

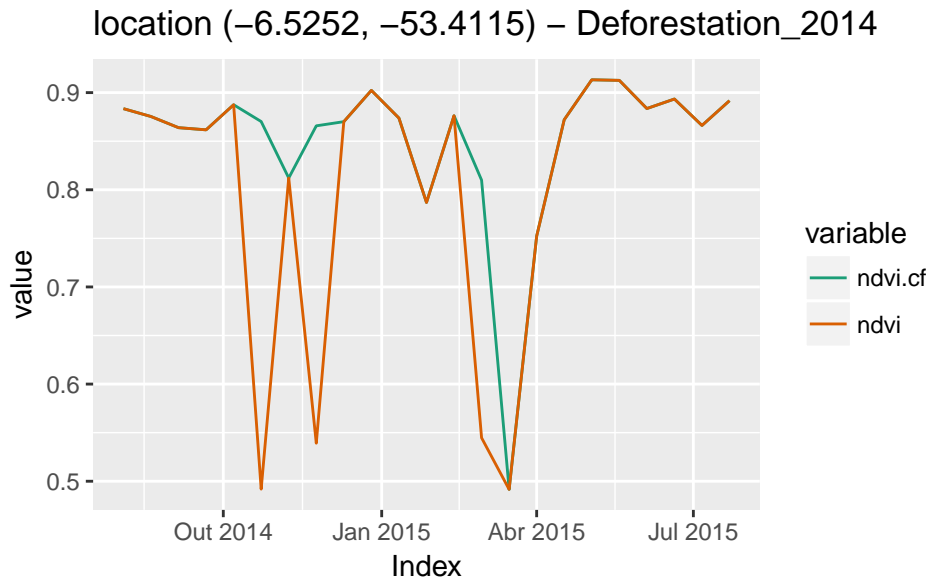
```
# Take the NDVI band of the first sample data set
point.tb <- sits_select(prodes_226_064[1,], bands = c("ndvi"))
# apply envelope filter (remove short downward noises)
point_env.tb <- sits_envelope(point.tb, "UL")
# plot the series
sits_plot(sits_merge(point_env.tb, point.tb))
```



Cloud filter

This function tries to remove noise from the input time series. It looks the first order difference time series for points where the value of the difference goes up abruptly. These points are taken as those whose difference is more than a cutoff value which is set by the user. Then, it applies an autoregressive integrated moving average (ARIMA) model to predict the missing values. The parameters of the ARIMA model can be set by the user. Please see `arima` for the detailed description of parameters `p`, `d`, and `q`.

```
# Take the NDVI band of the first sample data set
point.tb <- sits_select(prodes_226_064[1,], bands = c("ndvi"))
# apply ARIMA filter
point_cf.tb <- sits_cloud_filter(point.tb, apply_whit = FALSE)
# plot the series
sits_plot(sits_merge(point_cf.tb, point.tb))
```



Machine learning classification for land use and land cover using satellite image time series

The main advantage using satellite image time series in land use studies is that the time series is methodologically consistent with the very nature of the land covers. Using this kind of data allows focusing on land changes through time. Currently, most studies that use satellite image time series for land classification still use variations of the classical remote sensing image classification methods. Given a series of images, researchers use methods that produce a single composite for the whole series [Gomez et al., 2016]. In their review on this subject, Gomez et al. [2016] discuss 12 papers that use satellite image time series to derive image composites that are later used for classification. Câmara et al. [2016] denote these works as taking a *space-first, time-later* approach.

An example of *space-first, time-later* work on big EO data analysis is the work by Hansen et al. [2013]. Using more than 650,000 LANDSAT images and processing more than 140 billion pixels, the authors compared data from 2000 to 2010 to produce maps of global forest loss during the decade. A pixel-based classification algorithm was used to process each image to detect forest cover. The method classifies each 2D image one by one.

In our view, these methods do not use the full potential of satellite image time series. The benefits of remote sensing time series analysis arise when the temporal resolution of the big data set is able to capture the most important changes. Here, the temporal autocorrelation of the data can be stronger than the spatial autocorrelation. Given data with adequate repeatability, a pixel will be more related to its temporal neighbours than to its spatial ones. In this case, *time-first, space-later* methods lead to better results than the *space-first, time-later* approach [Câmara et al., 2016].

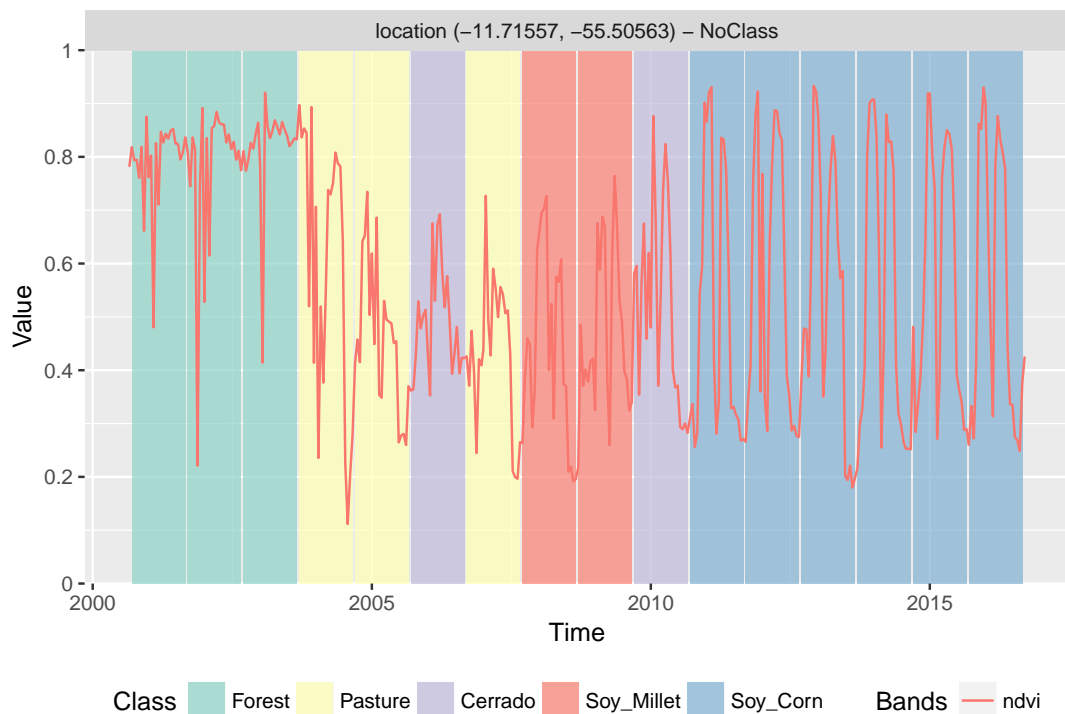
The `sits` package provides functionality to explore the full depth of satellite image time series data. It treat time series as a feature vector. To be consistent, the procedure aligns all time series from different years by its time proximity considering an given

cropping schedule. Once aligned, the feature vector is formed by all pixel “bands”. The idea is to have as many temporal attributes as possible, increasing the dimension of the classification space. In this scenario, statistical learning models are the natural candidates to deal with high-dimensional data: learning to distinguish all land cover and land use classes from trusted samples exemplars, also known as training data, to infer classes of a larger data set. In the next section we discuss about machine learning techniques supported in `sits` with more detail.

In the training stage, Additional “bands” can be computed to increase the distinction between classes. This method has a deceptive simplicity

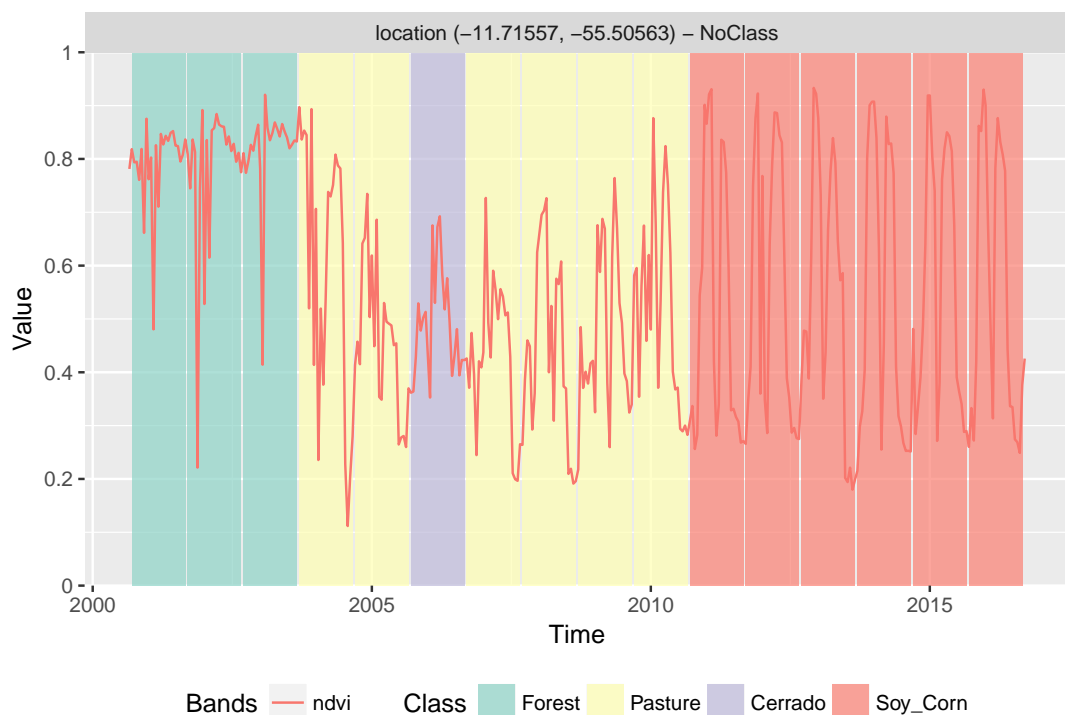
Support Vector Machine

```
# Retrieve the set of samples for the Mato Grosso region
# (provided by EMBRAPA) (samples_MT_ndvi) and
# get a point to be classified (point_ndvi)
class.tb <- sits_classify(point_ndvi,
                          samples_MT_ndvi,
                          ml_method = sits_svm(kernel = "radial",
                                                cost = 10))
sits_plot(class.tb)
```



Random Forest

```
# Retrieve the set of samples for the Mato Grosso region
# (provided by EMBRAPA) (samples_MT_ndvi) and
# get a point to be classified (point_ndvi)
class.tb <- sits_classify(point_ndvi,
                        samples_MT_ndvi,
                        ml_method = sits_rfor())
sits_plot(class.tb)
```



Validation techniques

Our experiment consists of the comparison of different methods to obtain the time series prototypes for each class. After obtaining the prototypes, we classified the data using the TWDTW method and used a cross-validation procedure to evaluate the results.

Final remarks

Current approaches to image time series analysis still use limited number of attributes. A common approach is deriving a small set of phenological parameters from vegetation indices, like beginning, peak, and length of growing season [Brown et al., 2013], [Kastens et al., 2017], [Estel et al., 2015], [Pelletier et al., 2016]. These phenological parameters are

then fed in specialised classifiers such as TIMESAT [Jönsson and Eklundh, 2004]. These approaches do not use the power of advanced statistical learning techniques to work on high-dimensional spaces and with big training data sets [James et al., 2013].

The sits uses the full depth of satellite image time series to create larger dimensional spaces. We tested different methods of extracting attributes from time series data, including those reported by Pelletier et al. [2016] and Kastens et al. [2017]. Our conclusion is that part of the information in raw time series is lost after filtering or statistical approximation. Thus, the method we developed has a deceptive simplicity: *use all the data available in the time series samples*. The idea is to have as many temporal attributes as possible, increasing the dimension of the classification space. Our experiments found out that modern statistical models such as support vector machines, and random forests perform better in high-dimensional spaces than in lower dimensional ones.

Acknowledgements

We want to thanks all the researchers whom provided data samples used in the examples. Alexandre Coutinho, Julio Esquerdo and Joao Antunes from the Brazilian Agricultural Research Agency (EMBRAPA) provided the samples for “soybean-fallow”, “fallow-cotton”, “soybean-cotton”, “soybean-corn”, “soybean-millet”, “soybean-sunflower” and “pasture” classes that were collected through farmer interviews and *in-loco* observations from October 2009 until October 2013. Rodrigo Bergotti from INPE whom provided samples for “cerrado” and “forest” classes through high resolution images observations. Damien Arvor [Arvor et al., 2012] whom provided ground samples for “soybean-fallow” class through *in-loco* observations.

References

- Saeed Aghabozorgi, Ali Seyed Shirkhorshidi, and Teh Ying Wah. Time-series clustering—a decade review. *Information Systems*, 53:16–38, 2015.
- D. Arvor, M. Meirelles, V. Dubreuil, A. Bégué, and Y. E. Shimabukuro. Analyzing the agricultural transition in Mato Grosso, Brazil, using satellite-derived indices. *Applied Geography*, 32(2):702–713, 2012.
- Peter M Atkinson, C Jeganathan, Jadu Dash, and Clement Atzberger. Inter-comparison of four models for smoothing satellite sensor time-series data to estimate vegetation phenology. *Remote sensing of environment*, 123:400–417, 2012.
- Clement Atzberger and Paul HC Eilers. Evaluating the effectiveness of smoothing algorithms in the absence of ground reference measurements. *International Journal of Remote Sensing*, 32(13):3689–3709, 2011.
- Bethany A Bradley, Robert W Jacob, John F Hermance, and John F Mustard. A curve fitting procedure to derive inter-annual phenologies from time series of noisy satellite ndvi data. *Remote Sensing of Environment*, 106(2):137–145, 2007.
- J. Christopher Brown, Jude H. Kastens, Alexandre Camargo Coutinho, Daniel de Castro Victoria, and Christopher R. Bishop. Classifying multiyear agricultural land use data from Mato Grosso using time-series MODIS vegetation index data. *Remote Sensing of Environment*, 130:39–50, 2013.
- Gilberto Câmara, Luiz Fernando Assis, Gilberto Ribeiro, Karine Reis Ferreira, Eduardo Llapa, and Lubia Vinhas. Big earth observation data analytics: matching requirements to system architectures. In *Proceedings of the 5th ACM SIGSPATIAL International Workshop on Analytics for Big Geospatial Data*, pages 1–6, Burlingame, CA, USA, 2016. ACM.
- Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):15, 2009.
- Jin Chen, Per. Jönsson, Masayuki Tamura, Zhihui Gu and Bunkei Matsushita, and Lars Eklundh. A simple method for reconstructing a high-quality ndvi time-series data set based on the savitzky–golay filter. *Remote Sensing of Environment*, 91(3-4):332 – 344, 2004.
- Johan de Rooi. *Penalized Estimation in High-Dimensional Data Analysis*. Ipskamp drukkers, Enschede, 2013.
- EMBRAPA. Sistema de Análise Temporal da Vegetação (SATVEG), 2014. URL www.satveg.cnptia.embrapa.br.
- Stephan Estel, Tobias Kuemmerle, Camilo Alcantara, Christian Levers, Alexander Prishchepov, and Patrick Hostert. Mapping farmland abandonment and recultivation across Europe using MODIS NDVI time series. *Remote Sensing of Environment*, 163: 312–325, 2015.

- Gillian L Galford, John F Mustard, Jerry Melillo, Aline Gendrin, Carlos C Cerri, and Carlos E Cerri. Wavelet analysis of MODIS time series to detect expansion and intensification of row-crop agriculture in Brazil. *Remote sensing of environment*, 112(2): 576–587, 2008.
- Cristina Gomez, Joanne C. White, and Michael A. Wulder. Optical remotely sensed time series data for land cover classification: A review. *{ISPRS} Journal of Photogrammetry and Remote Sensing*, 116:55 – 72, 2016. ISSN 0924-2716.
- M. C. Hansen, P. V. Potapov, R. Moore, M. Hancher, S. A. Turubanova, A. Tyukavina, D. Thau, S. V. Stehman, S. J. Goetz, T. R. Loveland, A. Kommareddy, A. Egorov, L. Chini, C. O. Justice, and J. R. G. Townshend. High-resolution global maps of 21st-century forest cover change. *Science*, 342(6160):850–853, 2013.
- T. Hastie, R. Tibshirani, and Friedman J. *The Elements of Statistical Learning. Data Mining, Inference, and Prediction*. Springer, New York, 2009.
- Christian Hennig. Clustering strategy and method selection. In Christian Hennig, Marina Meila, Fionn Murtagh, and Roberto Rocci, editors, *Handbook of cluster analysis*. CRC Press, 2015.
- Robert J. Hijmans. *raster: Geographic Data Analysis and Modeling*, 2015.
- Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 2(1):193–218, 1985.
- INPE. Amazon Deforestation Monitoring Project (PRODES), 2017. URL www.obt.inpe.br/prodes.
- G. James, D. Witten, T. Hastie, and R. Tibshirani. *An Introduction to Statistical Learning: with Applications in R*. Springer, New York, EUA, 2013.
- P. Jonsson and L. Eklundh. Seasonality extraction by function fitting to time-series of satellite sensor data. *IEEE Transactions on Geoscience and Remote Sensing*, 40(8): 1824–1832, 2002-08. ISSN 0196-2892. doi: 10.1109/TGRS.2002.802519.
- Per Jönsson and Lars Eklundh. TIMESAT—a program for analyzing time-series of satellite sensor data. *Computers & Geosciences*, 30(8):833 – 845, 2004. ISSN 0098-3004.
- J. Kastens, J. Brown, A. Coutinho, C. Bishop, and J. Esquerdo. Soy moratorium impacts on soybean and deforestation dynamics in Mato Grosso, Brazil. *PLOS ONE*, 12(4): e0176168, 2017.
- Robert E. Kennedy, Zhiqiang Yang, and Warren B. Cohen. Detecting trends in forest disturbance and recovery using yearly Landsat time series. *Remote Sensing of Environment*, 114(12):2897–2910, 2010.
- Eamonn Keogh, Jessica Lin, and Wagner Truppel. Clustering of time series subsequences is meaningless: Implications for previous and future research. In *Data Mining, 2003. ICDM 2003. Third IEEE International Conference on*, pages 115–122, 2003.
- E.F. Lambin and M. Linderman. Time series of remote sensing data for land change science. *IEEE Transactions on Geoscience and Remote Sensing*, 44(7):1926–1928, 2006.

- Eric F Lambin, Helmut J Geist, and Erika Lepers. Dynamics of land-use and land-cover change in tropical regions. *Annual review of environment and resources*, 28(1):205–241, 2003.
- T Warren Liao. Clustering of time series data—a survey. *Pattern recognition*, 38(11):1857–1874, 2005.
- Valerie J. Pasquarella, Christopher E. Holden, Les Kaufman, and Curtis E. Woodcock. From imagery to ecology: leveraging time series of all available LANDSAT observations to map and monitor ecosystem state and dynamics. *Remote Sensing in Ecology and Conservation*, 2(3):152–170, 2016. ISSN 2056-3485. doi: 10.1002/rse2.24.
- Charlotte Pelletier, Silvia Valero, Jordi Inglada, Nicolas Champion, and Gerard Dedieu. Assessing the robustness of random forests to map land cover with high resolution satellite image time series over large areas. *Remote Sensing of Environment*, 187:156–168, 2016.
- François Petitjean, Jordi Inglada, and Pierre Gançarskv. Clustering of satellite image time series under time warping. In *Analysis of Multi-temporal Remote Sensing Images (Multi-Temp)*, 2011 6th International Workshop on the, pages 69–72. IEEE, 2011.
- Gilberto Ribeiro de Queiroz, Karine Reis Ferreira, Lubia Vinhas, Gilberto Camara, Raphael Willian da Costa, Ricardo Cartaxo Modesto de Souza, Victor Wegner Maus, and Alber Sanchez. WTSS: um serviço web para extração de séries temporais de imagens de sensoriamento remoto. In *Proceeding of the XVII Remote Sensing Brazilian Symposium*, pages 7553–7560, 2015.
- Raoni Rajão and Niall Hayes. Conceptions of control and it artefacts: an institutional account of the amazon rainforest monitoring system. *Journal of Information Technology*, 24(4):320–331, 2009.
- LAS Romani, RRV Gonçalves, BF Amaral, DYT Chino, J Zullo, C Traina, EPM Sousa, and AJM Traina. Clustering analysis applied to ndvi/noaa multitemporal images to improve the monitoring process of sugarcane crops. In *Analysis of Multi-temporal Remote Sensing Images (Multi-Temp)*, 2011 6th International Workshop on the, pages 33–36. IEEE, 2011.
- Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- Toshihiro Sakamoto, Masayuki Yokozawa, Hitoshi Toritani, Michio Shibayama, Naoki Ishitsuka, and Hiroyuki Ohno. A crop phenology detection method using time-series MODIS data. *Remote Sensing of Environment*, 96(3-4):366–374, 2005.
- Abraham Savitzky and Marcel JE Golay. Smoothing and differentiation of data by simplified least squares procedures. *Analytical chemistry*, 36(8):1627–1639, 1964.
- Lubia Vinhas, Gilberto Ribeiro, Karine Reis Ferreira, and Gilberto Camara. Web services for big Earth observation data. In *Proceedings of the 17th Brazilian Symposium on GeoInformatics*, pages 26–35, Campos do Jordão, SP, Brazil, 2016. INPE.

- Joe H Ward. Hierarchical grouping to optimize an objective function. *Journal of the American statistical association*, 58(301):236–244, 1963.
- Edmund T Whittaker. On a new method of graduation. *Proceedings of the Edinburgh Mathematical Society*, 41:63–75, 1922.
- Hadley Wickham and Garrett Grolemund. *R for Data Science: Import, Tidy, Transform, Visualize, and Model Data*. O’Reilly Media, Inc., 2017.
- Achim Zeileis and Gabor Grothendieck. zoo: S3 infrastructure for regular and irregular time series. *Journal of Statistical Software*, 14(6):1–27, 2005.