

PIXELWISE TIME SERIES RETRIEVAL IN PHENOLOGICAL STUDIES

Elisangela Silva Santos¹, Bruna Alberton² and Leonor Patricia Morellato², and Ricardo da Silva Torres¹

¹Institute of Computing, University of Campinas, Campinas, Brazil

²Universidade Estadual Paulista, Rio Claro, SP, Brazil

ABSTRACT

The support of time series similarity searches might be crucial in phenology studies, in which long-term time series analysis based on the identification of similar and different phenological patterns shared by individuals belonging to different species is a widely common task. In this paper, we introduce the use of well-established Information Retrieval (IR) technologies in the search of time series. The solution comprises four main steps: extraction of an image-based time series representation; image content description to encode time series properties and patterns; textual signature extraction based on image content descriptions; and textual signature indexing using off-the-shelf IR approaches. In this paper, we demonstrate both the effectiveness and the efficiency of the proposed solution in time series retrieval problems related to the management of phenological data associated with near-surface vegetation images.

Index Terms— time series retrieval, recurrence plot, information retrieval, phenology

1. INTRODUCTION

Near-surface plant phenology stands for the use of a set of technologies (e.g., digital cameras) to capture plant phenological changes over time remotely [1]. Sequences of vegetation images are obtained aiming to support complex temporal change analyzes based on pixel properties, encoded in color channels. For example, plant leaf flushing properties are usually investigated by means of variations of the green channel over time, while leaf color change and senescence analyzes are often based on values from the red channel. In this context, searching for similar regions within images given their time series similarity is of paramount importance. Especially, in the tropics, where hundreds of species can be found within small-size regions, the identification of regions of interest, those associated within plant individuals whose temporal profiles have been kept track, plays an important role.

Authors are grateful to CNPq (grant #307560/2016-3), São Paulo Research Foundation – FAPESP (grants #2014/12236-1, #2015/24494-8, #2016/50250-1, and #2017/20945-0) and the FAPESP-Microsoft Virtual Institute (grants #2013/50155-0, #2013/50169-1, and #2014/50715-9). This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) - Finance Code 001.

The objective, in this case, is to identify specimens (or communities), whose phenological change patterns, encoded for example in time series associated with vegetation indices, are of interest.

This work is concerned with the definition of appropriate time series storage and retrieval mechanisms to support phenology studies. We address the time series similarity retrieval problem from both the effectiveness and efficiency perspectives. First, time series are represented using recurrence plot (RP) image representations [2], opening the opportunity for the use of effective image content description approaches for characterizing time series patterns. Second, we handle efficiency aspects by proposing to convert visual features into a textual signature [3], which is later indexed and searched using off-the-shelf information retrieval (IR) technologies. The use of RP is motivated by its recent success to encode complex non-linear patterns found within time series associated with phenological data [4].

Performed experiments considered plant phenology data related to near-surface vegetation images. The time series retrieval problem is defined in terms of regions within those images, which belong to the same plant species. Obtained results demonstrate that the proposed solution is effective and efficient for managing large volumes of times series, opening new opportunities for the investigation of IR technologies for time series indexing and retrieval. To the best of our knowledge, this is the first work to exploit the integration of RP representations, visual description approaches, and information indexing and retrieval technologies into a single searching tool.

2. TIME SERIES INDEXING/SEARCH SYSTEM

Figure 1 presents a schematic diagram of the proposed solution for indexing and searching time series. First, in (a) recurrence plot (RP) representations are computed. Next, regions and blocks within RP images are determined in (b). From each block, histogram-based visual features are extracted (c), and textual signatures are computed (d). Image textual signatures are concatenated (e) and indexed (f) for supporting searches (g). Next, these components are detailed.

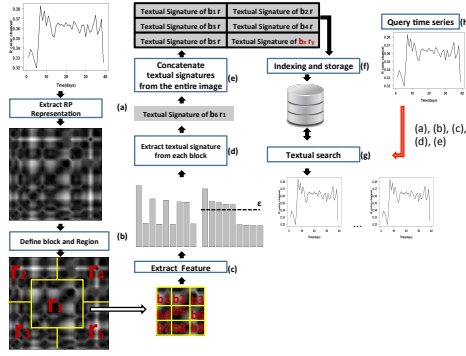


Fig. 1. Schematic diagram of the proposed solution.

2.1. Time Series Image Representation

Recurrence Plots (RP) [2] have been extensively used to characterize dynamical systems. Its use in nonlinear data analysis is motivated by its capacity of encoding repeated events of higher dimensional phase spaces. Typically, the *recurrence* of states are encoded into a two dimensional representation. This representation, referred to as $R_{i,j}$, can be defined as:

$$R_{i,j} = \Theta(\epsilon_i - \|x_i - x_j\|), x_i \in \mathcal{R}^m, i, j = 1 \dots N \quad (1)$$

where x is a time series, N is the time series size, ϵ_i is a threshold distance, $\|\cdot\|$ is a norm between the states (e.g., Euclidean distance), m is the embedding dimension, and $\Theta(\cdot)$ is the Heaviside function. This function has value 0 for negative argument and 1 for positive argument.

$R_{i,j}$ matrix can be seen as a square image (of dimension $N \times N$), usually know as *recurrence plot image*, which encodes complex temporal information. The created image might be binary or in gray scale, depending on the use (or not) of a threshold applied to the distance scores.

2.2. Image Content Description

The image content description follows the approach proposed in [3], which comprises two main steps: image partitioning, and histogram-based content characterization.

Let \hat{I}_{RP} be a recurrence plot image and $\mathcal{R} = \{r_1, r_2, \dots, r_\eta\}$ be a set of η partitions of \hat{I}_{RP} , such that each pixel p in \hat{I} belongs to only one region $r_j \in \mathcal{R}$. Let $\mathcal{B}_{r_j} = \{b_{1,r_j}, b_{2,r_j}, \dots, b_{\beta,r_j}\}$ be a set of non-overlapping blocks of consecutive pixels of region r_j . Figure 1 (bottom right) illustrates the definition of a 5-region partitioning within the RP image and nine blocks defined by a 3×3 grid within the central region.

Each block $b_{i,r_j} \in \mathcal{B}_{r_j}$ has its content represented by means of a histogram-based image description defined in terms of a function $D(b_{i,r_j})$, which assigns a feature vector $f_{v_{i,j}}$ to b_{i,r_j} . The use of color and texture-based histograms

were investigated in [3]. In this paper, we investigate the use of Histograms of Orientated Gradients (HOG) approach [5] and Local Binary Patterns (LBP) [6] for the characterization of region blocks.

2.3. Textual Signature Extraction

Each block $b_{i,r_j} \in \mathcal{B}_{r_j}$ has its content represented by means of a textual signature that is generated by a function $\delta(f_{v_{i,j}})$. This function is responsible for mapping a block feature vector $f_{v_{i,j}}$ to a visual word in the textual domain. The signature generated is defined as $\delta(f_{v_{i,j}}) = \langle \hat{r}_j - \hat{c}_1 - \hat{c}_2 - \dots - \hat{c}_n \rangle$, where $n \geq 0$, \hat{r}_j identifies the region r_j where the block b_i occurs and ‘-’ is a separator symbol. Each value \hat{c}_k encodes a content property (e.g., color, texture, or shape) that occurs in the block, such that the property frequency in the block is higher than a threshold ϵ and $(\hat{c}_i) > (\hat{c}_j)$, for $i < j$.

2.4. Signature Indexing

Textual signatures are indexed using traditional and widely used indexing schemes based on inverted files. In our implementation, we take advantage of freely available indexing libraries available in the Lucene system [7].¹

2.5. Time Series Search

Given a query time series, its textual signatures are extracted using the same procedures utilized for indexing (arrow highlighted in red in Figure 1). Ranked lists computation is based on the combination of a Vector Space Model (VSM) [8] and a Boolean model, available in the Lucene system.

3. EXPERIMENTS

This section presents performed experiments and discusses achieved results.

3.1. Experimental Protocol

3.1.1. Dataset

The dataset considered in this study refers to a sequence of vegetation images between August 29th and October 3rd 2011, day of year 241 to 278 (DOY), during the main leaf flushing season [9] of a Cerrado region located at Itirapina, São Paulo State, Brazil. These images were taken by a digital hemispherical lens camera (Mobotix Q24), installed on top of a 18m tower. The collection comprises daily sequence of five JPEG images (at 1280×960 pixels of resolution) per hour, from 6am to 6pm. We used the Guigues algorithm [10] to segment images hierarchically. The reference image used for segmentation was taken at noon on October 15th, 2011. Five segmentation scales from the hierarchy were considered

¹<http://lucene.apache.org/> (As of May 2019).

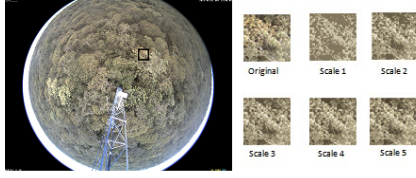


Fig. 2. Sample image recorded on Oct. 15th, 2011; and the segmentation results for the selected scales.

in order to extract time series. The finest scale, composed of 27,380 regions, were used for efficiency evaluation; while the coarsest scale, which contains 8,849 regions, was used for effectiveness assessment. Figure 2 illustrates the segmented scales in a subimage sample.

3.1.2. Time Series Computation

Time series are extracted from each region of the segmented image. Those time series refer to the variation of the contribution of the primary colors (Red, Green, and Blue) over time. As proposed by Richardson et al. [11], we compute, for each color channel, the average value of the pixel intensity. Later, the relative importance of each channel is computed as:

$$Total_{avg.} = Red_{avg.} + Green_{avg.} + Blue_{avg.} \quad (2)$$

$$Rcc = \frac{Red_{avg.}}{Total_{avg.}} \quad (3)$$

$$Gcc = \frac{Green_{avg.}}{Total_{avg.}} \quad (4)$$

$$Bcc = \frac{Blue_{avg.}}{Total_{avg.}} \quad (5)$$

3.1.3. Performance Assessment Criteria

Our strategy to evaluate both the effectiveness and the efficiency of the proposed search system relies on assessing the similarity among regions associated with individuals of the same species, based on their respective times series. We defined six ROIs based on the random selection of six plant species identified in the hemispheric image: (1) *Aspidosperma tomentosum*, (2) *Caryocar brasiliensis*, (3) *Myrcia guianensis*, (4) *Miconia rubiginosa*, (5) *Pouteria ramiflora*, and (6) *Pouteria torta*. We use, therefore, a retrieval protocol, in which, time series associated with regions belonging to individuals of the same species are relevant to each other. This protocol defines the ground truth for queries.

3.1.4. Evaluation Criteria

The effectiveness of the time series search system is computed in terms of the widely used mean Precision at 5 (Precision@5) for different queries. Precision@5 is defined as the rate of relevant time series retrieved in the top-5 positions of

Table 1. Effectiveness Results in terms of Precision@5 (Comparison with Baselines).

	Rcc	Gcc	Bcc
LBP	0.338 (7am)	0.400 (12am)	0.327 (6am)
HOG	0.290 (6am)	0.283 (6am)	0.329 (6am)
Fourier [12]	0.290 (7am)	0.249 (12am)	0.320 (6am)

the returned ranked list. The efficiency assessment is measured by the mean query processing time for answering a single time series query. We ran the experiments in a Intel Core i5-3230M CPU 2.60GHz \times 4 with 7.7 GB of memory, with Ubuntu OS.

3.2. Experimental Results

3.2.1. Effectiveness Evaluation

Figure 3 presents the effectiveness results in terms of Precision@5, considering the Rcc, Gcc, and the Bcc time series from 6am to 6pm. In all cases, usually, more effective results are observed early in the morning and late in the afternoon, regardless the descriptors used (either HOG or LBP). This is in accordance with previous results [12, 13]. One remarkable exception refers to the LBP effectiveness performance for 12am, which is associated with the highest score observed in all experiments (close to 40%).

Table 1 compares the effectiveness results of different configurations of the proposed time series retrieval system with a baseline. It presents the Precision@5 mean scores along with the time of the day related to the best performance. The baseline considered in this study is based on the use of the Fourier shape description approach for characterizing time series associated with vegetation indices. This method yielded the best results in the evaluation of shape description approaches for time series characterization presented in [12].

As we can observe, regardless the vegetation index, the proposed search scheme yielded the best results (highlighted in bold). LBP was most effective for Rcc and Gcc, while HOG was the best performer for the Bcc index.

3.2.2. Efficiency Evaluation

Figure 4 presents results related to efficiency aspects. This chart depicts the mean time in milliseconds to process a query. The performance of the search system increased linearly with the number of time series being indexed. The use of the LBP descriptor led to the most efficient results. Those results are even more impressive, given that a non-powerful machine was used in the conducted experiments.

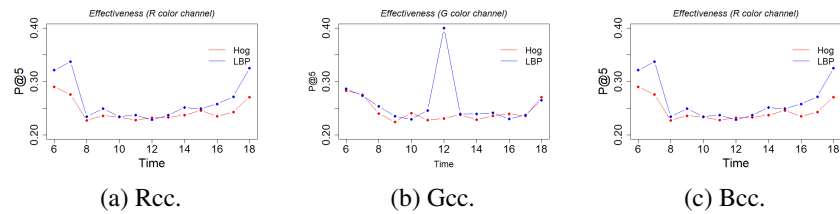


Fig. 3. Precision at 5 for (a) Rcc, (b) Gcc, and (c) Bcc vegetation indices.

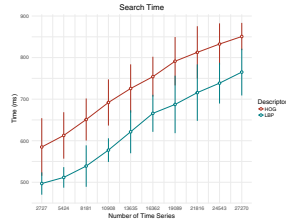


Fig. 4. Average response time for different collection sizes.

4. CONCLUSIONS

In this paper, we introduced a new approach for time series similarity searches. The proposed method relies on encoding time series into an image representation, whose content is characterized using image description approaches. Next, content properties are converted into textual signatures which can later be indexed and retrieved using state-of-the-art information retrieval approaches. Performed experiments involving time series related to plant phenology studies based on near-surface vegetation images demonstrate that the proposed method is effective and efficient in performing time series similarity searches. Future work includes the investigation of data-driven feature extractors in the characterization of recurrence plot image region blocks.

5. REFERENCES

- [1] B. Alberton, R. da S. Torres, L. F. Cancian, B. D. Borges, J. Almeida, G. C. Mariano, J. dos Santos, and L. P. C. Morellato, "Introducing digital cameras to monitor plant phenology in the tropics: applications for conservation," *Perspectives in Ecology and Conservation*, vol. 15, no. 2, pp. 82 – 90, Apr. 2017.
- [2] J.-P. Eckmann, S. Oliffson Kamphorst, and D. Ruelle, "Recurrence plots of dynamical systems," *Europhysics Letters*, vol. 4, no. 9, pp. 973, 1987.
- [3] J. M. dos Santos, E. S. de Moura, A. S. da Silva, and R. da S. Torres, "Color and texture applied to a signature-based bag of visual words method for image retrieval," *MTAP*, vol. 76, no. 15, pp. 16855–16872, Aug. 2017.
- [4] F. A. Faria, J. Almeida, B. Alberton, L. P. C. Morellato, A. Rocha, and R. da S. Torres, "Time series-based classifier fusion for fine-grained plant species recognition," *PRL*, vol. 81, pp. 101–109, October 2016.
- [5] Navneet Dalal and Bill Triggs, "Histograms of oriented gradients for human detection," in *CVPR*, 2005, pp. 886–893.
- [6] T. Ojala, M. Pietikainen, and T. Maenpaa, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE TPAMI*, vol. 24, no. 7, pp. 971–987, Jul 2002.
- [7] M. McCandless, E. Hatcher, and O. Gospodnetic, *Lucene in Action, Second Edition: Covers Apache Lucene 3.0*, Manning Publications Co., Greenwich, CT, USA, 2010.
- [8] Gerard Salton and M. J. McGill, *Introduction to Modern Information Retrieval*, McGraw-Hill, 1983.
- [9] B. Alberton, J. Almeida, R. Helm, R. da S. Torres, A. Menzel, and L. P. C. Morellato, "Using phenological cameras to track the green up in a cerrado savanna and its on-the-ground validation," *Ecological Informatics*, vol. 19, no. 0, pp. 62 – 70, January 2014.
- [10] L. Guigues, J. Cocquerez, and H. Le Men, "Scale-sets image analysis," *International Journal of Computer Vision*, vol. 68, pp. 289–317, 2006.
- [11] A. D. Richardson, J. P. Jenkins, B. H. Braswell, D. Y. Hollinger, S. V. Ollinger, and M. L. Smith, "Use of digital webcam images to track spring greening in a deciduous broadleaf forest," *Oecologia*, vol. 152, pp. 323–334, 2007.
- [12] R. da S. Torres, M. Hasegawa, S. Tabbone, J. Almeida, J. A. Santos, B. Alberton, and L. P. C. Morellato, "Shape-based time series analysis for remote phenology studies," in *IGARSS*, July 2013, pp. 3598–3601.
- [13] J. Almeida, J. A. dos Santos, W. O. Miranda, B. Alberton, L. P. C. Morellato, and R. da S. Torres, "Deriving vegetation indices for phenology analysis using genetic programming," *Ecological Informatics*, vol. 26, Part 3, pp. 61 – 69, 2015.