



PROGRAMA DE VERÃO DO LNCC

JORNADA DE CIÊNCIA DE DADOS

Apresentação do Docente



Eduardo Ogasawara

eogasawara@ieee.org

<https://eic.cefet-rj.br/~eogasawara>

Biografia



- Doutorado em Engenharia de Sistemas e Computação (COPPE/UFRJ) em 2011
- Docente do CEFET/RJ
 - Departamento de Ciência da Computação
 - Coordenação do Curso Técnico de Informática
- Docente permanente
 - Programa de Pós-graduação em Ciência da Computação (PPCIC)
 - Programa de Pós-graduação em Eng. de Produção e Sistemas (PPPRO)
- Membro da SBC, ACM, IEEE e INNS

<https://eic.cefet-rj.br/~eogasawara>

<http://lattes.cnpq.br/0528303491410251>

https://www.researchgate.net/profile/Eduardo_Ogasawara

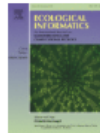
<https://www.linkedin.com/in/eogasawara>

Produção Científica: Centrada em Mineração de Dados e Workflows

- 125 artigos publicados*
 - 25 de periódicos



Ecological Informatics
Volume 36, November 2016, Pages 94-105





Evaluating temporal aggregation for predicting the sea surface temperature of the Atlantic Ocean



Knowledge-Based Systems
Volume 164, 15 January 2019, Pages 274-291



Nonstationary time series transformation methods: An experimental review

Rebecca Salles ^a, Kele Belloze ^a, Fabio Porto ^b, Pedro H. Gonzalez ^a, Eduardo Ogasawara ^a  

 Show more

<https://doi.org/10.1016/j.knosys.2018.10.041>

Get rights and content



Future Generation Computer Systems

Volume 68, March 2017, Pages 111-127



Deriving scientific workflows from algebraic experiment lines: A practical approach



Transportation Research Part E: Logistics and Transportation Review

Volume 95, November 2016, Pages 282-298



An analysis of Brazilian flight delays based on frequent patterns

Alice Sternberg ^a, Diego Carvalho ^a, Leonardo Murta ^c, Jorge Soares ^{a, b}, Eduardo Ogasawara ^a  

 Show more

<https://doi.org/10.1016/j.tre.2016.09.013>

Get rights and content

* Revisor de periódicos e conferências internacionais

Produção técnica: Diversos artefatos computacionais



Sim-Evolution

GPCA Educação



Hanafuda

GPCA Cartas



Este app é compatível com todos os seus dispositivos

ipeadata

BRASIL Acesso à Informação

macroeconômico regional social

Ministério da Ciência e Tecnologia

mineral data

Séries Históricas do Setor Mineral Brasileiro

Home Mapa do Site Fale conosco

Pesquisar: BUSCAR Acesso nº 66705.

Mineral Data
Substâncias
Temas
Fontes
Catálogo
Indicadores

Executor: Centro de Tecnologia Mineral - CETEM / Ministério da Ciência e Tecnologia - MC
Patrocinado pela Secretaria de Geologia, Mineração e Transformação Mineral - SGM / Ministério de Minas

STMotif: Discovery of Motifs in Spatial-Time Series

Allow to identify motifs in spatial-time series. A motif is a previously unknown subsequence of a (spatial) time series with relevant number of occurrences. For this purpose, the Combined Series Approach (CSA) is used.

Version: 1.0.4
Depends: R (≥ 2.10)
Imports: stats, ggplot2, reshape2, scales, grDevices, RColorBrewer, shiny
Suggests: knitr, rmarkdown, testthat
Published: 2019-08-22
Author: Heraldo Borges [aut, cre] (CEFET/RJ), Amin Bazaz [aut] (Polytech/Montpellier), Luciana Escobar [aut] (CEFET/RJ), Esther Pacitti [aut] (University of Montpellier), Eduardo Ogasawara [aut] (CEFET/RJ)

Maintainer: TSPred: Functions for Benchmarking Time Series Prediction

License: Functions for time series preprocessing, decomposition, prediction and accuracy assessment using automatic linear modelling. The generated linear models and its yielded prediction errors can be used for benchmarking other time series prediction methods and for creating a demand for the refinement of such methods. For this purpose, benchmark data from prediction competitions may be used.

Downloads: 40
Depends: R (≥ 2.10)
Reference: forecast, KFAS, stats, MuMin, EMD, wavelets, vars
Vignettes: 2018-06-21
Author: Rebecca Pontes Salles [aut, cre, cph] (CEFET/RJ), Eduardo Ogasawara [ths] (CEFET/RJ)

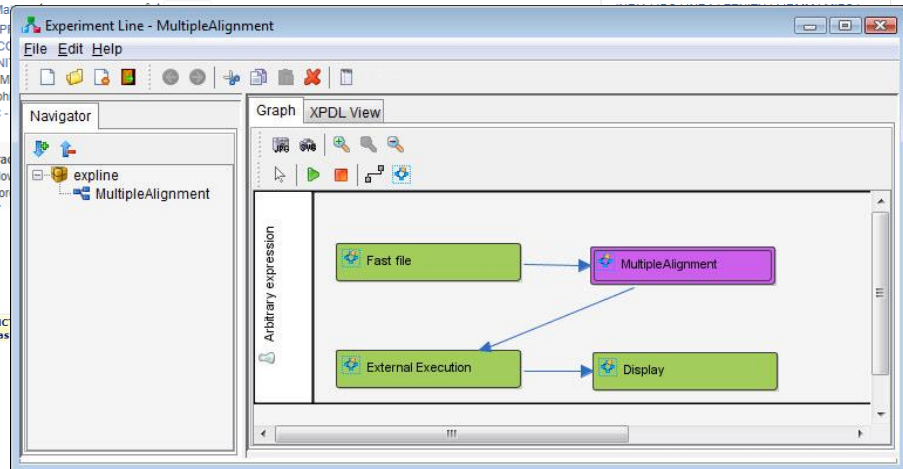
Home Official Web Site Document list Simple search

Package source: lirmm-00806557, version 1
Windows binaries: lirmm-00806557, version 1
OS X binaries: lirmm-00806557, version 1
Needs compilation: lirmm-00806557, version 1
Old source versions: lirmm-00806557, version 1

Chiron: A Parallel Engine for Algebraic Scientific Workflows

Download Eduardo Ogasawara¹, Dias Jonas¹, Vitor Silva¹, Chirigati Fernando¹, Oliveira Daniel De¹, Fabio Porto², Marta Marzari¹

Referer: COP
Package: LNCC
Window: ZEN
OS X binaries: LIRM
Old source: Soph
Old source: IBC-



Temas de Pesquisa: Predictive analytics

- Classificação
 - Regressão
 - Análise de séries temporais e espaço-temporais
 - Não-estacionariedade
- ## Descoberta e tratamento para *concept-drift*

On Evaluating Data Preprocessing Methods for Machine Learning Models for Flight Delays

Leonardo Moreira, Christofer Dantas, Leonardo Oliveira
CEFEFTR | leonardo.moreira,christofer.dantas,leonardo.oliveira@elc.cetec-tj.br

Jorge Soares, Eduardo Ogasawara
CEFEFTR | jorge@elc.cetec-tj.br | ogasawara@elc.cetec-tj.br

Abstract—Flight delays cause various inconveniences for airlines, airports, and passengers. According to data provided by the Brazilian National Civil Aviation Agency (ANAC), between 2009 and 2015, about 22% of domestic flights made in Brazil were delayed by more than 15 minutes. The prediction of these delays is fundamental to mitigate their occurrence and optimize the decision-making process of an air transport system. Particularly, airlines, airports, and users may be more interested in when delays are likely to occur than the accurate prediction of the absence of delays. This paper focuses on the unbalanced distribution of the classes of delay (presence and absence) by performing an experimental evaluation of several preprocessing methods for the development of machine-learning flight delay classification models. These models were built from a dataset that integrates national flight operations with meteorological conditions of airports. Our results indicate the models that applied the balancing techniques performed much better in predicting the occurrence of delays, getting about 60% of hits.

1. INTRODUCTION

Delay is one of the key performance indicators of any transportation system. A flight delay shall be represented by the difference between the programmed time and the actual time of departure or arrival of a flight. In the context of commercial aviation, these delays can occur for a variety of reasons, including flight process failures, weather conditions, mechanical problems, ground delays, air traffic control and capacity constraints.

In the commercial aviation scenario, delays have a high financial impact on airlines, such as fines, additional operating costs, and declining customer loyalty. Also, given the uncertainty of its occurrence, many passengers are forced to reschedule their travels to arrive at the destination on time, which often leads to increased travel costs [1].

Specialized literature shows that large volumes of data have been collected in databases of public and private institutions to study and to understand the operations of the air transport system. Analysis of this data is relevant for gaining the knowledge needed to detect and predict delays. Many recent studies have been done in analyzing flight data using machine learning methods [2], [3], [4]. Such initiatives are done thanks to the large volume of flight data that has been collected in these years, in what is currently known as the Data Science era [5], [6], [7], [8].

Thus, methods of predicting flight delays are fundamental to mitigate their occurrence, and, as a consequence, reduce the financial losses. Therefore, classification models for predicting

when delays may occur are needed, given the complexity of reasons and conditions that generate delays [9]. In fact, any improvement in this theme may be beneficial for airlines, airports, and passengers. Notably, the sensitivity of predicting when delays are likely to occur is more relevant than trying to target the accuracy of predicting the presence or absence of delays.

In this context, the objective of this work is to perform an experimental evaluation of data preprocessing methods for machine learning classification models, considering the factors involved and collected by the databases focusing on the sensitivity of machine learning classifiers. For that, this work builds a dataset that integrates a database containing flight operations data provided by Brazilian National Civil Aviation Agency (ANAC) [10] and airport weather data provided by Weather Underground [11]. From this dataset, many data preprocessing methods were applied in combination with different machine learning classification models. Their performance evaluation regarding delay prediction was analyzed.

This paper contributes by exploring a broader spectrum of data preprocessing methods for building machine learning models. Although flight delay prediction is an open problem, our results indicate that data preprocessing methods that target the problem of unbalanced distributions of the classes for delays outperform the other preprocessing methods. Additionally, we also contributed by exploring the Brazilian flight systems, where no other work was observed targeting predicting flight delays in this continental country [12].

Besides this introduction, this paper is organized as follows. Sections II and III present the general background for data preprocessing and machine learning methods. Section IV presents related work, whereas Section V discusses the methodology used for our exploratory analysis. Section VI analyzes the experimental evaluation. Finally, Section VII concludes and points out future work.

II. DATA PREPROCESSING

The data preprocessing activities allow for the input data to be prepared for the following data mining activities. Especially when it comes to classification, they also help in increasing both accuracy and performance of the machine learning methods. Among the main activities of data preprocessing we can highlight: Data Integration & Cleaning; (A) Data Transformation; (B) Data Reduction; (C) and (C)

Knowledge-Based Systems 164 (2019) 274–291

Contents lists available at ScienceDirect

Knowledge-Based Systems

Journal homepage: www.elsevier.com/locate/kbsys



Nonstationary time series transformation methods: An experimental review

Rebecca Salles¹, Kele Belloré¹, Fabio Porto², Pedro H. Gonzalez², Eduardo Ogasawara^{3,4}

¹Nacional Center for Technological Education of Rio de Janeiro (CEFET/RJ), Brazil

²Nacional Library for Scientific Computing (LNCC), Brazil

³CEFET/RJ, Brazil

⁴CEFEFTR, Brazil

⁵CEFEFTR, Brazil

⁶CEFEFTR, Brazil

⁷CEFEFTR, Brazil

⁸CEFEFTR, Brazil

⁹CEFEFTR, Brazil

¹⁰CEFEFTR, Brazil

¹¹CEFEFTR, Brazil

¹²CEFEFTR, Brazil

¹³CEFEFTR, Brazil

¹⁴CEFEFTR, Brazil

¹⁵CEFEFTR, Brazil

¹⁶CEFEFTR, Brazil

¹⁷CEFEFTR, Brazil

¹⁸CEFEFTR, Brazil

¹⁹CEFEFTR, Brazil

²⁰CEFEFTR, Brazil

²¹CEFEFTR, Brazil

²²CEFEFTR, Brazil

²³CEFEFTR, Brazil

²⁴CEFEFTR, Brazil

²⁵CEFEFTR, Brazil

²⁶CEFEFTR, Brazil

²⁷CEFEFTR, Brazil

²⁸CEFEFTR, Brazil

²⁹CEFEFTR, Brazil

³⁰CEFEFTR, Brazil

³¹CEFEFTR, Brazil

³²CEFEFTR, Brazil

³³CEFEFTR, Brazil

³⁴CEFEFTR, Brazil

³⁵CEFEFTR, Brazil

International Journal of Agricultural and Environmental Information Systems, 4(2), 23-36, April-June 2013 23

A Forecasting Method for Fertilizers Consumption in Brazil

Eduardo Ogasawara, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

Daniel de Oliveira, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

Fabio Paschoal Junior, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

Rafael Castaneda, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

Myrna Amorim, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

Renato Mauro, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

Jorge Soares, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

João Quadros, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

Eduardo Bezerra, CEFET-Federal Center of Technological Education Celso Suckow da Fonseca, Rio de Janeiro, Brazil

ABSTRACT

Tracking information about fertilizers consumption in the world is very important since they are used to produce agriculture commodities. Brazil consumes a large amount of fertilizers due to its large-scale agriculture fields. Most of these fertilizers are currently imported. The analysis of consumption of major fertilizers, such as Nitrogen-Phosphorus-Potassium (NPK), Sulphur, Phosphate Rock, Potash, and Nitrogen become critical for long-term government decisions. In this paper we present a method for fertilizers consumption forecasting based on both Autoregressive Integrated Moving Average (ARIMA) and logistic function models. Our method was used to forecast fertilizers consumption in Brazil for the next 10 years considering different economic growth for the entire country.

Keywords: Autoregressive Integrated Moving Average (ARIMA), Fertilizers Consumption, Forecast, Gross Domestic Product (GDP), Logistic Function, Population

DOI: 10.4018/ijaeis.2013040103

Copyright © 2013, IGI Global. Copying or distributing in print or electronic forms without written permission of IGI Global is prohibited.

Temas de Pesquisa: Workflows for data analytics

- Workflows
 - Modelagem: Atividades e dependências
 - Paralelização
- DSC (Data intensive scalable computing)
- Álgebra de workflow

An Algebraic Approach for Data-Centric Scientific Workflows

Eduardo Ogasawara^{1,2} Daniel de Oliveira¹ Patrick Valduriez²
Jonas Dias¹ Fábio Porto¹ Maria Mattoso¹
¹COPEPE/UFRRJ ²CEFET/RJ ³LNCC
Rio de Janeiro, Brazil Rio de Janeiro, Brazil Petrópolis, Brazil
ogasawara.jonasdias.daniel.marta@cos.ufrr.br fporto@lncc.br Patrick.Valduriez@lncc.br

ABSTRACT
Scientific workflows have emerged as a basic abstraction for structuring and executing scientific workflows in computational environments. In many instances, these workflows are computationally and data intensive, thus requiring execution in large-scale parallel computers. However, parallelization of scientific workflows remains low-level, ad-hoc and labor-intensive, which makes it hard to exploit optimization opportunities. To address this problem, we propose an algebraic approach (inspired by relational algebra) and a parallel execution model that enable automatic optimization of scientific workflows. We conducted a thorough validation of our approach using both a real oil exploration application and synthetic data scenarios. The experiments were run in Chaco, a data-centric scientific workflow engine implemented to support our algebraic approach. Our experiment demonstrates performance improvement of up to 25% compared to an ad-hoc workflow implementation.

1. INTRODUCTION

Many scientific experiments are based on complex computer simulations that consume and produce very large datasets and allocate huge amounts of computational resources. As the complexity of the experiments grows, running simulations becomes a challenge. To help scientists in managing resources involved in large-scale in-silico simulations, scientific workflows are gaining much interest. A workflow can be defined as a model of a process, which consists in a series of activities and its dependencies [1]. Workflows have been used primarily in business data processing. A data-centric scientific workflow, for short scientific workflow, structures the processing of a scientific simulation as a graph of activities, in which nodes correspond to data processing activities and edges represent the dataflow between them. Workflow activities are associated to scientific programs that prepare, process and analyze data.

Scientific Workflow Management Systems (SWMS) [2] are software systems that support the definition, execution and monitoring of scientific workflows. Various SWMSs have been

proposed (e.g. VisTrails, Kepler, Taverna, Pegasus, Swift and Triana). Each of them has its own language [3] and focuses on different aspects, such as parallel execution, semantic support, domain specific characteristics and management of provenance data.

Although some SWMSs focus on parallel execution, parallelizing large-scale simulations are still hard, ad-hoc and labor-intensive. Workflow developers (and scientists) need to decide on the ordering, dependencies, and the parallelization strategies. These decisions neglect parallelization opportunities, which may yield to miss important optimization opportunities. Let us illustrate the problem with a critical application we are addressing with researchers: Brian's game of the paper.

1.1 Motivating Example: RFA application
To illustrate the problem of optimizing data-centric scientific workflows, let us consider the following motivating workflow scenario from oil exploration. A major function of an ultra-deep water oil exploration system is pumping oil from thousands meters up to the surface through riser structures, called risers. Maintaining and repairing risers under deep water is difficult, costly and critical for the environment (e.g. to prevent oil spill). Understanding the dynamic behavior of each riser and its life expectancy is critical for Petrosbras. Thus, scientist must predict riser stages based on complex scientific models and observed data collected from risers. As shown in Figure 1, performing Riser Finite Element Analysis (RFA) requires a complex workflow of data-intensive activities that may take a very long time to complete. A typical riser's fatigue analysis workflow [3] consumes several input files containing riser information, such as finite element meshes, wind waves sea currents, case studies, and produces result analysis files to be further studied by the scientist. Such workflow can be modeled as a Directed Acyclic Graph (DAG) of activities with edges for establishing the dataflow. In Figure 1, each rectangle indicates an activity, solid lines represent input and output parameters that are considered between activities, and dashed lines represent input and output parameters shared through data files. For simplicity, we assume that all data is shared among all activities [4], with a shared disk storage. For a single riser's fatigue analysis, the workflow can have as many as 2,000 input files (about 1GB) and produces about 4,500 files (about 2GBs) with nine activities including:

¹ Work partially sponsored by CAPES, CNPq and DORA (Debrago and Sauer's project).

Rumo à Integração da Álgebra de Workflows com o Processamento de Consulta Relacional*

Jólio Antonio Ferreira¹, Jorge Soares¹, Fábio Porto¹,
Esther Paçótti¹, Rafael Coutinho¹, Eduardo Ogasawara¹

¹CEFET/RJ - Centro Federal de Educação Tecnológica Celso Suckow da Fonseca

²LNCC - Laboratório Nacional de Computação Científica

³lnria - University of Montpellier

joao.pazara@lum.org, jorge@lncc.br, rafael@lncc.br, fporto@lncc.br,

Esther.Paco@lnria.fr, rafas@lncc.cottinno@cefet-rj.br, ogasawara@ieee.org

Abstract. *Workflows emerged as a basic abstraction for structuring data analysis experiments in the current Data Intensive Scalable Computing (DISC) scenario. In many situations, these workflows are intensive, either computationally or in relation to data management, requiring execution in high-performance processing environments. However, parallelizing the execution of workflows commonly requires laborious programming, in an ad hoc manner and in a low level of abstraction, which makes it difficult to explore optimization opportunities. Some algebraic approaches have been developed to mitigate such limitation. This work moves in the direction converging the workflow algebra with relational query processing.*

Resumo. *Os workflows emergiram como uma abstração básica para estruturar experimentos de análise de dados no atual cenário de DSC (Data Intensive Scalable Computing). Em muitas situações, esses workflows são intensivos, seja computacionalmente ou em relação à manipulação de dados, exigindo a execução em ambientes de processamento de alto desempenho. Entretanto, paralelizar a execução de workflows comumente requer programação trabalhosa, de modo ad hoc e em baixo nível de abstração, o que torna difícil a exploração das oportunidades de otimização. Algumas abordagens algébricas foram desenvolvidas visando mitigar tal limitação. Esse trabalho caminha na direção de convergir a álgebra de workflows com o processamento de consultas relacionais.*

1. Contexto

Apesar de alguns sistemas de workflows possuírem recursos para execução paralela, paralelizar um workflow de larga escala é uma tarefa difícil, ad hoc e trabalhosa. Na maioria das soluções existentes, cabe aos usuários dos sistemas decidirem a ordem e as dependências entre as atividades além das estratégias de paralelização. Estas decisões, em muitos casos, resultam em oportunidades de otimização da execução do workflow que poderiam levar a melhorias significativas de desempenho [Ogasawara et al., 2011]. Principalmente

*Os autores agradam a FAPERJ, FCAJES e ao CNPq pelo financiamento parcial do projeto.

August 27-26, 2016 - Rio de Janeiro, RJ, Brazil

Exploring Machine Learning Methods for the Star/Galaxy Separation Problem

Eduardo Machado
Marcelo Senjeira
Dairaldo Ogasawara
CEFEPR/UFRRJ
ogasawara@iaee.org

Ricardo Ogando
Márcio A. G. Maia
Lairi Nicolaci da Costa
Observatório Nacional, LInEA
{ogando, lmaoia, maia}@lnln.gov.br

Ricardo Campisano
Gustavo Paiva Guedes
Eduardo Bezerra
CEFEPR/UFRRJ
{guedes, ebezerra}@cefepr-ufrrj.br

Abstract. *For recent or planned deep astronomical surveys, it is important to tell stars and galaxies apart, a task known as Star/Galaxy Separation Problem (SGSP). At faint magnitudes, the separation between pointy and extended sources is fuzzy, which makes SGSP a hard task. This problem is even harder for large surveys like Dark Energy Survey (DES) and, in a near future, the Large Synoptic Survey Telescope (LSST) due to their large data volume. Hence, the search for classification methods that are both accurate and efficient is highly relevant. In this work, we present a comparative analysis of several machine learning methods targeted at solving the SGSP at faint magnitudes. In order to train the classification models, the COSMOS survey was used. We use machine learning methods as distinct as artificial neural networks, k-nearest-neighbor, Support Vector Machines, Random Forests and Naive Bayes. The exploratory process was modeled as data-centric workflow. The workflow was implemented on top of Hadoop framework and was used to find the best parameter values for each classification method we considered, of which neural networks and random forest present superior performance.*

1. INTRODUCTION

A lot has changed since the first astronomical surveys [1] when, along with the well-known stars and planets, a new class of extended sources, the nebulae, was discovered. It was only in the beginning of the XX century that this new class was later associated to extra-galactic sources, the galaxies. This distinction between extended and point sources was enough to tell them apart until recently.

As astronomical surveys pushed the boundaries of the observed universe, reaching fainter magnitude limits and challenging the spatial resolution of ground-based telescopes, this distinction between stars and galaxies became fuzzier, as one just can not spatially resolve distant galaxies, confining them with stars or quasars.

One of the main issues that affect the spatial resolution of ground-based telescopes is the blur caused by atmosphere turbulence, or seeing. Removing the atmosphere from the line of sight using space telescopes, like Hubble Space Telescope (HST) [2], one can resolve distant galaxies that would be otherwise indistinguishable from stars when observed from Earth.

The discovery of the accelerated expansion of the universe [3] caused by an unknown component generally called Dark Energy, started a race to uncover its nature. Several large and

deep photometric surveys are ongoing (DES [4]) or under construction (LSST [5]) with the main goal of studying dark energy. In order to do that, cosmological probes such as large scale structure, weak lensing, and galaxy clusters, require a very good distinction between stars and galaxies. While surveys like SDSS [6] relied successfully on morphological classification on its data releases, there is an urge for new techniques when dealing with deeper surveys. Many works explore machine learning techniques [7], [8], specially using cheap HST images from COSMOS survey [2] as training set. As the amount of data collected in deeper surveys grows, the demand for fast and automatic star/galaxy separation methods is increasing. Many studies exist that use one or two classification methods. In particular, artificial neural networks have been a common tool [7], and, more recently, support vector machines [8]. On the other hand, it is well-known that no algorithm performs well on all possible classification problems [9]. To the best of our knowledge, the literature for star/galaxy separation problem lacks (1) a more extensive experimental analysis of a broader spectrum of machine learning methods, specially near the faint/unsolvable end of magnitude/shape distribution, and (2) a disciplined approach to fine-tune the parameters of such methods.

In this paper, we explore the application of several machine learning methods to the star/galaxy separation problem, namely, neural networks, SVM, random forests, k-nearest neighbors and Naive Bayes. To train and validate these methods, we use COSMOS survey [2], which consists of a catalog of more than 50,000 objects and 90 attributes. The exploratory process was modeled as data-centric workflow and implemented on top of Hadoop framework. Our approach takes into account preprocessing activities and parameter exploration that are not commonly used in this scenario. Our experiments show that neural networks and random forest present superior performance than the remaining methods. The result is important because random forest is not among the most used methods in star/galaxy separation problem, which reinforces the need for such exploration.

This work is organized in the following sections. Section II presents an overview on Astronomy surveys and related concepts. Section III provides a brief description of the processing techniques whereas Section IV presents the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted by ACM provided that the fee code for users to reproduce and distribute for profit or commercial advertising and for resale is paid to ACM and that the base fee code on this page for copying exceeds the fee code for users to reproduce. In prior work we have used the same methodology as in this paper, but we did not use specific permissions and we are thankful to the reviewers who pointed to the need of the ACM International Conference on Very Large Data Bases, Proceedings of the VLDB Endowment, Vol. 4, No. 12, Copyright 2010 VLDB Endowment 2159-5886/10 \$10.00.

Temas de Pesquisa: Pattern mining

- Mineração de seqüências
- Descoberta de motivos
- Detecção de eventos, anomalias, outliers



An analysis of Brazilian flight delays based on frequent patterns
Alice Sternberg^a, Diego Carvalho^a, Leonardo Murta^c, Jorge Soares^{ab}, Eduardo Ogasawara^{a,b*}

ARTICLE INFO
Received 23 March 2016
Received in revised form 31 August 2016
Accepted 30 September 2016

ABSTRACT
In this paper we applied data indexing techniques combined with association rules to unveil hidden patterns of flight delays. Considering Brazilian flight data and guided by six research questions related to causes, moments, differences, and relationships between airports and airlines, we evaluated and quantified all attributes that may lead to delays, showing not only the main patterns, but also their chances of occurrence in the entire network. In each airport and airline, we observed that Brazilian flight system has difficulties to recover from previous delays and when operating under adverse meteorological conditions, delays occurrences may increase up to 216%.

1. Introduction
Delays are one of the greatest challenges to transportation systems. Notably, in commercial aviation, delay is usually defined by the difference between scheduled and real times of departure or arrival (Wieland, 1997). Despite some differences in tolerance thresholds for delays, country regulatory authorities usually monitor delays through several indicators. In 2014, 16.6% of flights delayed by more than 15 min in Europe and 24.7% in the United States. In Brazil, 19.1% of domestic flights were canceled or suffered delays greater than 30 min (EUROCONTROL, 2015; The United States Department of Transportation, 2015; ANAC, 2015b).
Flight delays impact passengers, airlines, and airports, especially increasing trip and operations costs. Given the uncertainty of their occurrence, passengers usually plan to travel earlier to ensure their arrival on time. On the other hand, airlines may have to pay penalties, fines, or incur extra cost, such as crew reschedules and aircraft relocations in airports (Britto et al., 2012). Moreover, delays are also related to environmental damages, since they may increase fuel consumption and gas emissions (Pejovic et al., 2009; Ryerson et al., 2014; Simić and Bahić, 2015).
Delays also affect the airlines marketing strategies, since the loyalty of customers are motivated by punctual performances (Vlachos and Lin, 2014). Furthermore, delay levels are not only related to operational and economic choices of an airline (such as aircraft sizes, flight frequencies, and fares), but also with complaints about airline service (Bhadra, 2009; Pai, 2010; Zou and Hansen, 2014). In this context, understanding the reasons for flight delays occurrences can direct public and private investments in air transportation systems, improve tactical and operational decisions of airports and airlines managers, and warn passengers, so they can rearrange their plans (Marsden, 2002; Lv and Wang, 2009).
Every moment, a massive amount of data from commercial aviation is collected and stored in public and private databases. Seeking to understand the air transportation ecosystem, domain analysts and data scientists are intensifying the usage

* Corresponding author.
E-mail address: ogasawara@ieee.org (E. Ogasawara).

Discovering Tight Space-Time Sequences

Ricardo Campisano¹, Heraldo Borges¹, Fabio Porto², Fabio Perosi³, Esther Pacitti⁴, Florent Massegla⁴, and Eduardo Ogasawara^{1(✉)}

¹ CEPET/RJ - Federal Center for Technological Education of Rio de Janeiro, Rio de Janeiro, Brazil
² LNCC - National Laboratory of Scientific Computing, Petrópolis, Brazil
³ UFRJ - Federal University of Rio de Janeiro, Rio de Janeiro, Brazil
⁴ Inria and University of Montpellier, Montpellier, France

Abstract. The problem of discovering spatiotemporal sequential patterns affects a broad range of applications. Many initiatives find sequences constrained by space and time. This paper addresses an appealing new challenge for this domain: find tight space-time sequences, i.e., find within the same process: (i) frequent sequences constrained in space and time that may not be frequent in the entire dataset and (ii) the time interval and space range where these sequences are frequent. The discovery of such patterns along with their constraints may lead to extract valuable knowledge that can remain hidden using traditional methods since their support is extremely low over the entire dataset. We introduce a new *Spatio-Temporal Sequence Miner (STSM)* algorithm to discover tight space-time sequences. We evaluate STSM using a proof of concept use case. When compared with general spatial-time sequence mining algorithms (*GSTSM*), STSM allows for new insights by detecting maximal space-time areas where each pattern is frequent. To the best of our knowledge, this is the first solution to tackle the problem of identifying tight space-time sequences.

1 Introduction

Space and time are pervasive in our day-to-day lives. As many datasets that include both time and space data are becoming available, new opportunities to discover interesting spatiotemporal patterns arise. An event may be classified as an occurrence of a phenomenon in a given space and time. A spatiotemporal sequential pattern is a sequence of events that are constrained in space and time [7]. Due to that, spatiotemporal sequence mining is gaining attention [11, 12].
In this work, we investigate a new problem related to spatiotemporal pattern identification. We are interested in finding tight space-time sequences, i.e., sequences that are constrained in space and time that may not be frequent in the entire dataset but are frequent inside a time interval and space range (spatiotemporal blocks). The primary challenge is to discover these blocks and the frequent sequences they contain. Solving this problem has a valuable impact on many applications.



Detecção de Anomalias Frequentes no Transporte Rodoviário Urbano¹

Ana Beatriz Cruz¹, João Ferreira¹, Diego Carvalho¹, Eduardo Mendes¹, Esther Pacitti², Rafael Coutinho³, Fabio Porto², Eduardo Ogasawara¹

¹ CEFET/RJ
² LNCC - DEXL Lab
³ FGV
⁴ Inria & University of Montpellier

anacruz@acm.org, joao.parana@acm.org, d.carvalho@ieee.org,
eduardo.mendes@fgv.br, renato.souza@fgv.br, Esther.Pacitti@inria.fr,
rafael11.coutinho@cefet-rj.br, fport@lncc.br, ogasawara@ieee.org

Abstract. The growth of urban population and, consequently, the number of vehicles causes the increase of traffic jams and emission of polluting gases. In this context, we observe the intensification of papers that aim to identify bottlenecks and their causes. These papers propose methodologies that use trajectory data model and aim to explain systemic behaviors. This article proposes the identification and classification of anomalies in the urban road transport system from space-time aggregations to permanent objects. The methodology consists of pre-processing of data, identification of anomalies, identification, and classification of frequent patterns. Through it, we can identify the systemic and specific behaviors on the urban transit of Rio de Janeiro.

Resumo. O crescimento da população urbana e, consequentemente, do número de veículos provoca o aumento de engarrafamentos e da emissão de gases poluentes. Nesse contexto, observamos a intensificação de pesquisas que buscam identificar engarrafamentos e suas causas. Essas pesquisas propõem metodologias que usam modelo de dados de trajetória e visam explicar comportamentos sistêmicos. Este artigo propõe a identificação e a classificação de anomalias no sistema de transporte rodoviário urbano a partir de agregações espaço-temporais a objetos permanentes. A metodologia consiste do pré-processamento dos dados, identificação de anomalias, identificação e classificação de padrões frequentes. Por meio dela, é possível identificar comportamentos sistêmicos e pontuais do trânsito urbano do Rio de Janeiro.

1. Introdução

Em 2007, pela primeira vez existiam mais pessoas vivendo em áreas urbanas do que em zonas rurais, resultado de uma urbanização expressiva que se impulsionou desde a década

¹ Os autores agradecem à FAPERJ, à CAPES e ao CNPq pelo financiamento parcial do projeto.

Colaboração

- Institutos de Pesquisa
 - LNCC
 - Fabio Porto, Artur Ziviani, Antônio Tadeu A. Gomes
 - Fiocruz
 - Marcel Pedroso, Cristiano Boccolini, Christovam Barcellos
- Academia
 - CEFET/RJ
 - A maioria dos docentes do PPCIC/PPPPO
 - COPPE/UFRJ
 - Marta Mattoso, Geraldo Zimbrão, Geraldo Xexéo
 - UFF
 - Daniel Oliveira, Leonardo Murta, Vanessa Braganholo
- Internacionais
 - INRIA / University of Montpellier
 - Patrick Valduriez, Esther Pacitti, Florent Masegla, Reza Akbarinia