

FACULDADE DE ENGENHARIA DA UNIVERSIDADE DO PORTO

Social Media Text Processing and Semantic Analysis for Smart Cities

João Filipe Figueiredo Pereira



**FACULDADE DE ENGENHARIA
UNIVERSIDADE DO PORTO**

Mestrado Integrado em Engenharia Informática e Computação

Supervisor: Rosaldo José Fernandes Rossetti

Supervisor: Pedro dos Santos Saleiro da Cruz

June 15, 2017

Social Media Text Processing and Semantic Analysis for Smart Cities

João Filipe Figueiredo Pereira

Mestrado Integrado em Engenharia Informática e Computação

Abstract

With the rise of Social Media, people obtain and share information almost instantly on a 24/7 basis. Many research areas have tried to extract valuable insights from these large volumes of freely available user generated content. The research areas of intelligent transportation systems and smart cities are no exception. However, extracting meaningful and actionable knowledge from user generated content is a complex endeavor. First, each social media service has its own data collection specificities and constraints, second the volume of messages/posts produced can be overwhelming for automatic processing and mining, and last but not the least, social media texts are usually short, informal, with a lot of abbreviations, jargon, slang and idioms.

In this thesis, we try to tackle some of the aforementioned challenges with the goal of extracting knowledge from social media streams that might be useful in the context of intelligent transportation systems and smart cities. We designed and developed a framework for collection, processing and mining of geo-located Tweets. More specifically, it provides functionalities for parallel collection of geo-located tweets from multiple pre-defined bounding boxes (cities or regions), including filtering of non complying tweets, text pre-processing for Portuguese and English language, topic modeling, and transportation-specific text classifiers, as well as, aggregation and data visualization.

We performed empirical studies and implemented illustrative examples for 5 cities: Rio de Janeiro, São Paulo, New York City, London and Melbourne, comprising a total of more than X millions tweets in a period of 3 months. The topic modeling and text classifiers were evaluated with manual labeled data specifically created for this work. Both software and gold standard data will be made publicly available to foster further developments from the research community.

Resumo

Devido à ascensão das Redes Sociais, as pessoas obtêm e partilham informação quase que instantaneamente 24/7. Muitas áreas de investigação tentaram extrair informações importantes destes grandes volumes de conteúdo, gerado por utilizadores, e livremente disponíveis. As áreas de investigação de sistemas inteligentes de transportes e de cidades inteligentes (*smart cities*) não são exceção. Contudo, extrair conhecimento acionável e significativo de conteúdo gerado por utilizadores exige um esforço complexo. Primeiro, cada serviço de social media possui as suas próprias especificidades e restrições para o método de recolha dos dados; em segundo lugar, o volume de mensagens produzidas pode ser esmagador para o processamento automático e prospeção; e por último, não menos importante, os textos das redes sociais são, geralmente, curtos, informais, com muitas abreviações, jargões, gírias e expressões idiomáticas.

Nesta dissertação, tentamos abordar alguns dos desafios acima mencionados com o objectivo de extrair conhecimento de mensagens das redes sociais que possam ser úteis no contexto de sistemas inteligentes de transportes e cidades inteligentes (*smart cities*). Nós idealizamos e desenvolvemos uma *framework* para a recolha de dados, processamento e prospeção de Tweets geo-localizados. Mais especificamente, a *framework* fornece funcionalidades para a recolha paralela de tweets geo-localizados de *bounding-boxes* (cidades ou regiões), incluindo filtragem de tweets não preenchidos, pré-processamento de texto para a língua portuguesa e inglesa, modelagem de tópicos e classificadores de texto específicos para transportes, bem como, agregação e visualização de dados.

Realizamos estudos empíricos e implementamos exemplos ilustrativos para 5 cidades: Rio de Janeiro, São Paulo, Nova York, Londres e Melbourne, perfazendo um total de mais de X milhões de tweets em um período de 3 meses. O modelo de tópicos e os classificadores de texto foram avaliados com dados manualmente anotados e criados especificamente para este trabalho. Tanto os dados quanto o software criados serão disponibilizados publicamente para promover novos desenvolvimentos da comunidade de investigação.

Acknowledgements

João Pereira

*“You should be glad that bridge fell down.
I was planning to build thirteen more to that same design”*

Isambard Kingdom Brunel

Contents

CONTENTS

List of Figures

LIST OF FIGURES

List of Tables

LIST OF TABLES

Abbreviations

SC	Smart City
SM	Smart Mobility
ITS	Intelligent Transportation System
ICT	Information and Communication Technology
SMA	Social Media Analytics
HTTP	Hypertext Transfer Protocol
TSL	Transport Security Layer
POS	Part-of-speech
BoW	Bag-of-words
VSM	Vector Space Model
LDA	Latent Dirichlet Allocation
CRF	Conditional Random Fields
HHM	Hidden Markov Model
ABSA	Aspect-based Sentiment Analysis
SSWE	Sentiment Specific Word Embeddings
ML	Machine Learning
SVM	Support Vector Machines
NB	Naïve Bayes
ME	Maximum Entropy
RF	Random Forests
DL	Deep Learning
MAE	Mean Absolute Error
OLS	Ordinary Least Squares
LR	Logistic Regression

Chapter 1

² Introduction

⁴ 1.1 Context and Motivation

The rise of social media services, in the last few years, has led to an excessive amount of information being placed online by the population. The need to explore this type of information has steadily grown in order to realize what kind of value could bring to the areas of marketing, business or even politician [?]. Micro-blogging platforms, such as Twitter, have a huge affluence in a daily-basis, where people publicly share around 500 million messages about a diverse set of current themes, expressing their opinions and feelings [?]. For this reason, technological projects developed in some cities have seen social media streams as a potential resource to extract knowledge, i.e. the satisfaction of people regarding some services in the cities, such as the urban transportation service [?]. The collection of this participative activity of people through online platforms, named Crowd Sensing, emerged to replace the traditional way of sensing data capture and tracking in physical infrastructures, allowing a considerable reduction of economic costs [?].

¹⁶ The information extraction from social media streams is a hard task. For example, tweets besides being text messages and, consequently unstructured data. There's also some extra particularities, as, for example, the limited length (140 characters per message) which restricts the amount of information in its content, the informal language used and there are also many spelling mistakes, abbreviations, ambiguity and special mentions (e.g. URL references, hashtags, users) and the presence of a high variety of emoticons [?].

²² Many research projects have been conducted in order to extract the sentiment present in opinions through text mining techniques. Sentiment analysis is the field that focuses on this task. The detection of the sentiment polarity in messages generated by citizens, whether at a specific level (relative to certain aspects) or at a general level (general sentiment of the message) can allow companies or even ordinary citizens the possibility of identifying city's services problems and may be assimilated to a kind of sensor for the issue of quality awareness, providing, at least, some help in decision-making processes.

1.2 Problem Statement and Goals

30 The problem around this dissertation establishes in the analysis of a continuous flow of social
media streams, in particular from Twitter, about a target scenario, as, for example, the quality of the
32 urban transportation services in Porto. Hence, it will be necessary to filter the relevant messages,
extract sentiment aspects/topics with polarity and create useful aggregated data visualizations. The
34 problem presented can be divided in five distinct points:

1. Data collection for our target scenario

36 At this point, a real scenario must be chosen so that a case study can be produced.

2. Named Entity Disambiguation and content filtering

38 The identification of the entities present in the messages is an important point in order to dis-
ambiguate some mentions that could not be related with target scenario, making the filtering
40 task easier, since only messages that are related must appear in the final dataset.

3. Identification of aspects/topics in Twitter messages

42 Since each opinion usually has a target aspect/topic about an entity, or even a service of
a city, the recognition of it is relevant so that the sentiment present in the message has an
44 orientation.

4. Sentiment polarity classification

46 The polarity of a sentiment may have three different types: positive, negative or neutral. At
this point, it will be necessary to estimate the level of polarity expressed in the message.

5. Data aggregation and visualization

50 The aggregation of the results provided by all other tasks is needed. Some messages could
refer to the same aspect through different ways, so it will be necessary aggregate messages
52 that present this characteristic. Another important task of the aggregation is the continu-
ously calculation of the results So that when an user access the analytics UI, the results
54 are presented immediately, without waiting time. At the visualization of the results, some
qualitative and quantitative indicators may be presented to the end-user to make the analysis
easier.

56 Taking into account all the aforementioned points, the final goal of this dissertation is to create
a framework, based on the concept of analysis. The framework should be capable of automati-
58 cally processing social media texts, regarding semantic processing, topic detection and sentiment
analysis. An user interface will be provided to the end-user, illustrating a set of qualitative and
60 quantitative indicators to analysis. This knowledge can be relevant both to users of a particular
service or to the responsible entities in order to improve the decision-making process.

62 1.3 Structure of this Dissertation

64 This report covers a great diversity of points and because of that, its structure is divided in three
different sections.

66 The Section ?? starts with a brief contextualization in the Smart Cities and Intelligent Trans-
portation System fields. After that, it was made a review on Social Media Analytics, specially
about Twitter, and what benefits its exploration can provide to the civilization. To explore the
68 information from Twitter messages, an intensive study regarding the Text Mining area was made,
in particular Information Extraction, Topic Modeling and the various Sentiment Analysis fields,
70 as well as related works that already explored similar problems with social media data.

72 The proposed solution for the aforementioned problem and its methodology are referenced in
Section ???. The final section of this planning report is composed by some conclusions relatively
to the studied works and what benefits and risks our solution may have - Section ??.

Introduction

⁷⁴ **Chapter 2**

Background and Literature Review

⁷⁶ This section aims to analyse and reflect about some works and topics that will be relevant to fully understand the problem. The study of solutions found by other authors can simplify the difficult
⁷⁸ task that is the analysis of social media data. Hence, this section has been divided into several parts in order to perceive not only the environment in which the problem is located but also the
⁸⁰ most important points to be studied in order to build our the final product. Respectively to the problem scope its important to know what is a smart city and how the transportation system can
⁸² contribute to this meaning. Since our product is a framework which goal is to extract information about social media data, i.e. texts from the Twitter, its interesting obtain the knowledge about
⁸⁴ extraction tools, such as the Twitter APIs, in order to have an idea how to construct our crawler module. The meaning of text mining and how the information present in texts can be extracted
⁸⁶ through different kinds of techniques, regarding disambiguation, filtering and modeling. Finally, it is also important to analyze several works about sentiment analysis in order to know the different
⁸⁸ methodologies and which are the most advantageous for this problem.

2.1 Smart Cities and Intelligent Transportation Systems

⁹⁰ The Smart City concept appeared thanks to the continuous growth of a city's population which has contributed to an aggressive urbanization [?]. In the last few years, several definitions for its
⁹² meaning have emerged, but the ideal one is not yet fully known. Angelidou in [?] defines Smart City as a "conceptual urban development model on the basis of the utilization of human, collective,
⁹⁴ and technological capital for the development of urban agglomerations" and enhance as its primary key the knowledge and the innovation economy. In her work, there is an identification of four
⁹⁶ forces that model the concept of a Smart City and two of them are very important to enhance: *technology push*, where new products and solutions are introduced in the market regarding the fast
⁹⁸ advance in science and technology; *demand pull*, where solutions and problems are developed in order to respond to the society demands, like the continuous growth of the population [?].

¹⁰⁰ The development environment in a city tagged with the concept "smart" is another key factor to reach the success. Komninos focus the importance of collective sources of innovation to the

Background and Literature Review

102 improvement of life quality in cities. The globalization of innovation networks are the responsible
103 for the emergency of another types of environments, such as "global "innovation clusters and
104 i-hubs, intelligent agglomerations, intelligent technology districts and intelligent clusters, living
105 labs" making possible the experimentation of products or services by the population in order to
106 identify problems or even to analyse the behavior of the people regarding what they have experimen-
107 tered [?].

108 The transportation system is inherently connected to the progress of a city, since people on a daily-basis uses the several ways of transportation, i.e. bus, private cars, metropolitan, etc, to go to
109 their jobs and make their own life. This system is also influenced by the problem of the population
110 growth being relevant the need of finding solutions to minimize or even raze it [?]. Hence, "a
111 smart city should be focused on its actions to become smart", coming up the concept of innovation
112 [?].

113 To understand what are *Intelligent Transportation Systems*, it is crucial introduce the meaning
114 of Smart Mobility. SM is a combination of comprehensive and smarter traffic service with
115 smart technology, enabling several intelligent traffic systems which provide control in the signals
116 regarding the traffic volume, information about smooth traffic flows, times of bus, train, subway
117 and flight arrivals and their routes [?]. The majority of *Intelligent Transportation Systems* are ex-
118 pressed through smart applications where the transportation and traffic management has became
119 more efficient and practicable, allowing the users to access important information about the trans-
120 portation systems in order to make correct decisions about what they want to use in their cities
121 [?]. ICT-based infrastructures are the main support for Smart Cities when the focus are ITS, since
122 through that is possible to pilot the activities operations and its management over a long period of
123 time [?].

124 Nowadays, cities are exploring some initiatives of sensing to support the development of tech-
125 nological projects. Areas such as utilities management (where, for example, is monitored the
126 consumption level of power, water and gas), traffic management (using vibration sensors to mea-
127 sure the traffic flows on bridges, or even the full capacity of a parking lot), environment awareness
128 (using video cameras to monitor the population behaviour and sensors to measure the level of air
129 pollution) make use of physical sensors, i.e. some devices that can capture information to study
130 and improve the quality of life in a daily basis [?]. R. Szabo et al. [?] and D. Doran et al. [?]
131 report the highly economic cost that this kind of sensing needs, since it's require the maintenance
132 and replacement of this devices, as well as a tracking infrastructure store and treat the information
133 collected. Hence, a new form of sensing has emerged - Crowd Sensing - to offer the cities several
134 ways to improve their services by exploring the participation of the population in the social net-
135 works where there are a publicly share of citizen's opinions and thoughts regarding some problems
136 [?]. This type of sensing consists in *human-generated* data provided by the population through the
137 use of mobile devices and the social networks platforms. Such data can be further used to extract
138 some analytics regarding specific services in a city, namely the urban transportation system [?].
139 Based on all this, social media can be seen as a good source of data to extract valuable information
140 in order to direct it to the smartness evolution process of a city [?].

142 Several works have already been developed and presented taking into account these two large
143 areas, Transportation Services and Smart Cities, using social media as source of information. G.
144 Anastasi et al. [?] proposed a framework which objective was the promotion of flexible transporta-
145 tion systems usage, i.e. encouraging people to share transport or to opt for the use of bicycles in
146 order to minimize infrastructural and environmental problems. Their tool takes advantages of the
147 crowd sensing techniques by exploring social media streams to predict accidents or traffic conges-
148 tion and alert the users of their service about this type of events. W. Liu et al. [?] have made a
149 study in three different transportation modes (private cars, public transportsations and bicyclists)
150 using theirs channels on Twitter to estimate a percentage of the majority gender that uses this
151 services in the city of Toronto. They have extract all the channel's tweets appealing only to the
152 *non-protected* followers and applied an already developed classification model to label each tweet
with its creator gender: male or female.

153 T. Ludwig et al. [?] proposed a tool capable of collect and display social media streams in order
154 to help the integration and coordination of volunteers in actions performed by emergency services
155 to prevent engagement in dangerous areas. Their tool present to the end-users map visualization
156 of a city where they could identify public calls of the emergency services to accept or deny them.

157 In conclusion to everything that has been analyzed in this sub-section, it's possible to verify
158 that the cities are increasingly opting for technological opportunities that involve crowd sensing,
159 once this type of exploration brings a considerable reduction of costs and the information that is
160 collected may contribute to the extraction of value from data generated by the population itself.

161 **2.2 Social Media Analytics**

162 In the last few years social networks have made impact on the business communications, since
163 users assume the role of costumers through the publication of content on this networks, rising the
164 levels of interaction between users and businesses entities [?]. A proof of that is the amount of
165 information produced since 2011 which is equivalent to a number over than 90% of the available
166 data online [?]. Facebook, Twitter and other social networking sites are nowadays used as business
167 tools by companies aiming the efficient use of digital marketing techniques to publicize their products
168 [?]. Besides the business field, the population turn into this new communication technologies
169 in a intensely way, where they publicy share real-life events, their opinions about certain topic,
170 their on-time feelings in the network through a simple message [?]. Social Media Analytics can
171 be describe as a type of digital analytics to study the people interaction with others, or their opin-
172 ion about companies, its products and services through the social media data. This study provides
173 important information to "analysts, brands, agencies or vendors", and its analysis could facilitate
174 the generation of economic value to many organizations [?]. To achieve the main goals of the
175 SMA, the companies focus their effort in the development and evaluation of frameworks, to make
176 possible an easy collection, analyse, summarization and visualization of processed social media
177 data. Hence, the companies can establish specific points about what to improve in their products
178 [?]. To create a significantly value regarding the SMA, J. Philips in [?] enhance some important

180 factors: users permissions, the listening of real-time information, the search mechanism, the data
access and integration, and others, before the choice of a tool that allows the information collec-
182 tion. Besides the tool, is also important have an idea of what is need to explore because the use
of a wrong technique of SMA could have bad business impact for the company. The majority of
184 SMA techniques focus on modeling in order to understand the large range of data collected and
support techniques, such as sentiment analysis, trend analysis and topic modeling, are the most
186 commonly used [?].

2.2.1 Twitter

188 Twitter is a social network where people freely micro-blogging about any topic and, like any
other social network, makes possible the connection between users around the world [?]. This
190 social network has faced an exponential growth since its inception, and nowadays its users, which
surpass 200 millions, produce around 500 millions tweets daily, performing a massive bunch of
192 information that could be an ideal testbed for research projects on big data [?, ?]. N. Banerjee
et al. [?] classify Twitter as a micro-blogging service presenting three attributes that justify such
194 characterization.

- **Limited Context Information:** The length of a Twitter message never exceeds 140 charac-
196 ters, which, in cases of knowledge extraction, since the amount of information is short, the
final results could not be the expected and none knowledge contribute is obtained.
- **Richness of Exchange:** Symbolizes the great diversity of posts that exists in the micro-blog
198 network. People talk about their daily activities, have conversations with each other and
shares their thoughts (moods or opinions) about a certain event or topic.
- **High Dimensionality:** The informality and ambiguity present in the messages and the ex-
202 panse vocabulary present in any language makes a large dimension of data. The informality
of the messages can be seen as the presence of spelling errors and abbreviations in its con-
tent, while the ambiguity can be the presence of words with multiple meanings.

Twitter has not only evolved in terms of usability, but also in the purpose of its use, i.e. Twitter
206 is not only used as personal diary for people but also represents a source of information for re-
porters and journalists to find potential news about real-time events [?]. A good example is given
208 by J. Sankaranarayanan et al. [?] in the period of Michael Jackson's death. The first tweet about
the incident was posted 20 minutes after the call to the 911 emergency service and, nearly, two
210 hours before the first communication on the news as it's possible to verify in the figure ??.

One of the advantages of Twitter compared to other social networks, for example, Facebook,
212 is the easier way to access its users originated data. While Facebook does not provide private
information about its users unless there are permissions to do so, or the content shared is present
214 in public pages or groups, Twitter allows the collection of all tweets from channels or directly
from people in order to be analysed in any kind of project because the user's accounts are usually
216 public [?, ?].

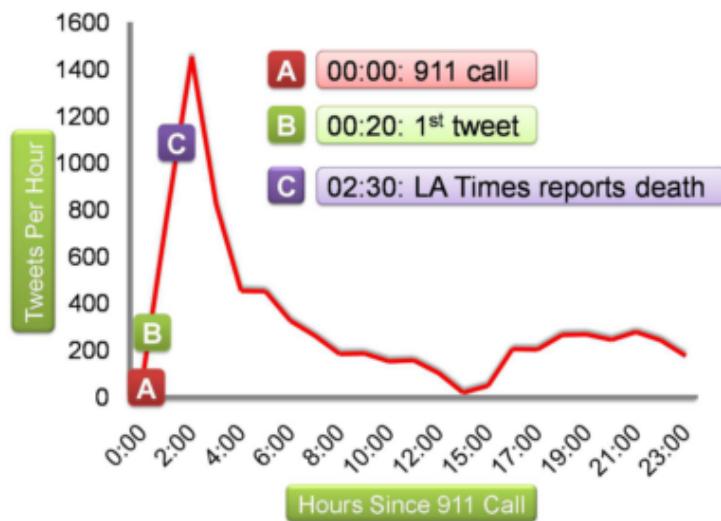


Figure 2.1: Traffic in tweets per hour relating to Michael Jackson's death by J. Sankaranarayanan [?]

Although this freedom in the collection of data, Twitter has also an ethical perspective and a
 218 regulation must be accomplish by developers or researchers. The TOS (Twitter Terms of Service)
 was created in order to make known what can be done with the data and protect the users' rights
 220 [?].

2.3 Text Mining

222 Text mining is a derived field from Data mining and aims to extract valuable information from
 unstructured textual data[?]. The reason why this technology is nowadays so much explored is
 224 because of the massive amount of information that is stored in text documents, such as "text
 files, HTML files, chat messages and emails" and it's required an automated technique that make
 226 possible the identification, extraction, management, integration and the knowledge exploration of
 information from texts in a efficiently and systematically way [?]. On the other hand, the social
 228 media applications also have contributed to the growth of text mining usage where companies have
 seen a potential path to improve their business model and increase the economic value relatively
 230 its competitors.

A. Stavrianou et al. [?] identify text mining as an interdisciplinary field since this technology
 232 takes advantages from Data mining techniques and combines several methodologies from simi-
 lar research areas, such as Categorization, Information Extraction, Information Retrieval, Topic
 234 Tracking and Concept Linkage. A common problem related to text mining is its similarity with
 Information Retrieval and Information Extraction which leads people to a non-differentiation be-
 236 tween this technologies. The difference between Information Retrieval and Text mining is estab-
 lished in their final goal, while IR aims to find and retrieve documents that match a certain part of

Table 2.1: Text Mining Issues by A. Stavrianou [?]

Issue	Details
Stop list	Should we take into account stop words?
Stemming	Should we reduce the words to their stems?
Noisy Data	Should the text be clear of noisy data?
Word Sense Disambiguation	Should we clarify the meaning of words in a text?
Tagging	What about data annotation and/or part of speech characteristics?
Collocations	What about compound or technical terms?
Grammar / Syntax	Should we make a syntactic or grammatical analysis? What about data dependency, anaphoric problems or scope ambiguity?
Tokenization	Should we tokenize by words or phrases and if so, how?
Text Representation	Which terms are important? Words or phrases? Nouns or adjectives? Which text model should we use? What about word order, context, and background knowledge?
Automated Learning	Should we use categorization? Which similarity measures should be applied?

238 a text or some keywords (e.g. Google Search Engine¹), TM tries to discovery unknown patterns
in texts that can be interpreted and explain some facts or truths contained in the lexical [?, ?, ?].

240 Regarding the Information Extraction, the differentiation can be seen in the data specificity and
structure. IE focus on the extraction of expected information from structured data and precocious
242 relations, while the information returned by TM techniques should be unsuspected and unexpected
with the data holding an unstructured format [?].

244 The motivation behind text mining holds on the benefit that other fields of research could take
from a use of its techniques. Information Retrieval systems can improve their precision since its
246 basis is the identification of semantic relations. Several areas can explore this methodology to find
inconsistencies in relational databases and make the integration, update and querying tasks easier
248 [?].

Text mining shares some of the issues presented by the Natural Language Processing field.
250 Once texts are usually performed by humans some associated problems can appear, such as
spelling mistakes, wrong phrasal construction, slang among other. Before the "mining" of a text,
252 it's important to apply some pre-processing steps in order eliminating noisy data from the primary
analysis process. A. Stavrianou et al. cite this issues very well in they work and it can be seen in
254 Table ??.

The removal of words from the text can sometimes not be desirable because some sentences
256 can lose its information or even leads to a different meaning compared with its original form. The

¹<https://google.com>

generation of a stop list words should be a supervised task as long as little words could induce
258 distinct results in the text classification [?].

Stemming is a task that depends mostly from the language of the text than its domain [?] and
260 the main goal of this technique is to reduce a word to its root to help in the calculus of distances
between texts or even keywords or phrases.

262 The noisy data is derived from spelling mistakes, acronyms and abbreviations in texts and to
solve this, a conversion of this terms should be done to keep a valid integrity of the data. The most
264 commonly solution approaches involve text edit distances (Levenshtein Distance²) and phonetic
distances measures between known words and the misspelling ones to achieve good corrections [?]

266 Word Sense Disambiguation focus on solving the meaning ambiguity present in words. Other
similar field to WSD is Name Entity Disambiguation (NED) where the disambiguation target are
268 named-entities mentions, while WSD focus on common words. WordNet³ is a resource very used
to extinguish this ambiguity [?]. There are two types of disambiguation, the supervised, where the
270 task is support by a dictionary or a thesaurus [?], and the unsupervised one, where the different
meanings of a word are unknown and normally learning algorithms with training examples are
272 used to achieve good results in the disambiguation task [?].

274 Tagging can be describe as the process of labeling each term of the text with a part-of-speech
tag, i.e. classify each word as a noun, verb, adjective, etc [?]. Collacation are very important in
text mining, since this task consist in group two or more words to give the correct meaning content
276 in the text. Collacations are usually made before the WSD task since some compound technical
terms have different meaning from the individual words which composed it [?].

278 Tokenization serves to pick up all the terms presented in a text document and to achieve this it's
necessary the split of the document into a stream of words implying the removal of the punctuation
280 marks and non-text characters [?]. some authors also see tokenization as a text representation form
since one of the most used models to represent texts is *Bag-of-words* (BoW). This model broke
282 down texts into words and stores it in a vector being also presented the word frequency occurrence
in the text. Hence, each word may represent a feature [?]. Another commonly used model to rep-
284 resent texts is Vector Space Models that represent all the documents in a multi-dimensional space
where documents are converted to vectors and each vector may be seen as a feature. This model
286 provides some advantages since the documents can be compared with each other by performing
some specific vector operations [?].

288 The purpose of this section is to provide the reader the definition of what is text mining, and
also the identification of basic operations and steps that are necessary for the preprocessing of this
290 type of unstructured data - texts.

²https://en.wikipedia.org/wiki/Levenshtein_distance

³<https://wordnet.princeton.edu/>

2.4 Information Extraction

292 Information Extraction is an important field of text mining and its main goal is finding structured
 293 information from semi-structured or unstructured texts. This kind of information can range from
 294 the identification of entities, such as people, organizations and places names, to a relation between
 295 this concepts. In the sentence, "*At 1976, Apple was founded by Steve Jobs and his friends*", its
 296 possible extract information about who were the founders of Apple and what was the year of its
 297 foundation. Problems regarding the analysis of the sentence can be located on the words "Apple"
 298 and "his", and how could a machine know that "Apple" is a reference of a technological company
 299 and not a reference to a fruit, or even, that the word "his" establish a connection between "Steve
 300 Jobs" and "friends". The exploration of this kind of information constitutes a potential measure to
 301 improve computer systems, such as search engines and database management [?].

302 When the target analysis is social media data, i.e. texts from social networks, the filtering
 303 process of the data is a crucial step. It's desirable that only related-topic information should be
 304 collected in order to avoid the existence of noisy objects in the data set [?, ?].

Information Extraction presents a set of components that can tackle this problems, such as
 305 Named Entity Recognition (NER), Named Entity Disambiguation (NED) and among others.

2.4.1 Name Entity Recognition and Name Entity Disambiguation

308 NER and NED are two distinct tasks that sometimes could cause some ambiguity in its purpose.

Named Entity Recognition is seen as a sub-task of Information Extraction aiming the correct
 310 labeling of words in a text in order to have knowledge of its types. The accuracy retrieve by this
 311 task is important when further steps depends on it, e.g. Relation Extraction [?]. Gazetteers are
 312 commonly used in this task since they provide pre-processed lists of organizations, days, places
 313 and person names which can be matching against some terms we want to recognize. In some
 314 cases, this "tools" are not enough to solve the problem because their domain is very limited and
 315 some terms can not exist, implying the use of external knowledge to fill this lack [?].

316 There are two main approaches to conduct this task: Ruled-based and Statistical Learning. The
 317 Ruled-based approach focus on the definition of a set of features for each token in the text to further
 318 comparing the text with a bunch of rules. The rules are composed by patterns that should trigger
 319 some labeling action in a sequence of tokens. The definition of this rules usually requires human
 320 expertise [?]. The Statistical Learning approach, also known as Statistical Machine Learning,
 321 treats the text as a sequence of observations which are represented by a vector of features. The
 322 final goal of this approach is the assignment of a label y_i to each observation x_i . The mapping
 323 process usually follows a BIO notation, firstly introduced to text chunking by [?], where the entity
 324 name could be at the beginning (B) or inside (I) of the observation and never outside (O). Patawar
 et al. [?] enhance three types of methods to this approach.

326 Supervised methods, where it's characteristic the existence of a labeled training data set to train
 327 a model and then classify a set of test to measure the model performance and accuracy. Hidden
 328 Markov Model was used in [?] to recognize and classify some text. In [?], Decision Trees were

combined with a simple rule generator to prove that this method could achieve similar results as
 330 Maximum-Entropy-based methods used in [?]. Support Vector Machines were used by H. Isozaki
 et al. [?] where they have proved that SVM models could also achieve good results for NER tasks
 332 if some analysis was carefully made in the *kernel functions* and filtering methods were applied,
 like the removal of useless features.

334 Semi-supervised methods used different amount of training data, i.e. labeled examples, and
 test data. The test data is usually a bigger amount facing the training data. The methodology
 336 commonly applied to this approach involves the use of "bootstrapping" which is an iterative pro-
 cess of training the model with progressive supervised increases of the training data set until the
 338 performance starts to decrease [?].

The last type are the Unsupervised Methods, where a large amount of labeled data is necessary
 340 and it's difficult to have this requirement. This need bases on the huge number of features required
 to this kind of methods. To fill this lack, a very frequent method used is the clustering where the
 342 formation of groups is made using the similarity present in the texts domain [?].

Y. Li et al. [?] define Named Entity Disambiguation as a "process of associating an entity
 344 name mentioned in a text to an entry, representing that entity, in a knowledge base". In the last
 few years, NED has been a target for a considerable number of research projects. The majority
 346 methods implemented to tackle this task are focused in three main features. Y. Li et al. [?] enhance
entity popularity as a statistical one where there is a "assumption that the most prominent entity
 348 for a given entity mention is the most probable underlying entity for that mention". At this feature,
 a link between the term and its Wikipedia page reference is established. The *context similarity* is
 350 another feature and aims the complementarity of the *entity popularity* feature. This feature centers
 on similarity measures between the text in analysis and the content text of the Wikipedia page. Y.
 352 Li reveals that this feature is word-dependency since it's necessary that both texts shares identical
 words in order to produce expected results. *Topical Coherence* is the third feature and solves the
 354 emerged problem of the second one. This feature uses the Wikipedia cross-page links mechanism
 in order to look up for related-topics of other entities and makes a connection with the target entity
 356 in the disambiguation problem. Through this process the domain text is expanded, decreasing the
 word-dependency problem appeared in the second feature.

358 D. Spina et al. [?] present two different approaches to solve the problem of ambiguity pre-
 sented in texts. The first one is *entity linking* and consists in the establishment of an association
 360 between the mention in the text and a entity present in a Knowledge Base. Three steps are needed
 to perform the linking of the entity name to the knowledge base:

- 362 • **Query Expansion:** Mining the Wikipedia structure and solving co-references in the docu-
 ment in order to enrich the query;
- 364 • **Candidate Generation:** Construction of a list of candidate entities from the knowledge
 base according to the information presented in the query;
- 366 • **Candidate Ranking:** It is the phase in which is computed some similarity measures be-
 tween the query and the entities, in order to rank the candidates and select the best one.

368 Another approach to tackle this problem of disambiguation of entities presented by D. Spina
et al. [?] is the *document enrichment by linking to Wikipedia Articles*. Similar to the *entity linking*,
370 this approach is also composed by three different steps and takes advantages from text representa-
tion models, such as Bag-of-words or Vector Space Models, for example.

- 372 • **Mention or surface form representation:** There a context definition of the target mention
to disambiguate. Text representation models are build with all set of entities that are am-
374 biguous and the unambiguous ones which are already resolved with linking to Wikipedia
pages.
- 376 • **Candidates entities retrieval and representation:** All the candidate entities referenced by
the knowledge base (e.g. Wikipedia or Freebase) and the content on the page are converted
378 to text representation models. After this, there is extraction of some features that can range
from page categories of the candidate entities or even syntactic features.
- 380 • **Best candidate selection:** The computing of similarity and distance functions between the
two text representation models, produced in the two previous steps, is made to select the
382 best candidate.

Some works related to NED were made in the context of micro-blogging services, such as
384 Twitter.

At 2010, Ferragina et al. [?] developed a system capable of identity entities in short texts
386 as, for example, tweets. Their system take advantage of the hyperlink mechanism of Wikipedia,
extracting related links between pages and the anchors texts in the links. By detecting some senses
388 present in the anchors, they try to disambiguate the ambiguous ones through a collective agreement
function, i.e. a voting classification. They used the unambiguous senses to boost the selection of
390 the ambiguous ones and have trying some pruning in the anchors set to improve the performance
of the system.

392 Meji et al. [?], similar to Ferragina et al., also explored anchors texts in Wikipedia articles.
The authors have used a supervised machine learning approach to conceive a list of candidates to
394 disambiguate each mention present in their tweets. Their strategy focus on the identification of
some patterns in each tweet, such as n-grams, to further matching it with the anchor texts of the
396 Wikipedia articles, taking also in account the hyperlink mechanism of this Knowledge Base.

Considering this works it's possible conclude that Wikipedia is a potential source to explore
398 in order to solve mentions of entities that could lead to a ambiguous problem.

2.4.2 Content Filtering

400 The content filtering is one of the most important tasks to analyze micro-blog data (e.g. Twitter).
The main goal of content filtering is the classification of Twitter posts which contains an entity
402 name, assuming the existence of a relation between the name and the content in order to erase
ambiguity in the dataset [?]. Recently, some contests related with Online Reputation Monitoring

Background and Literature Review

404 (ORM) have explored this task of filtering content. The WePS-3 ⁴ and the RepLab 2012 tackle
405 unknown-entity scenario approaches, while the RepLab 2013 ⁵ focus on the known-entity scenario
406 approach.

In the WePS-3, the LSIR [?] research group has build a system where a profile identify each
408 of the companies mention. The Wordnet ⁶ and the company web-page were used to extract a
410 bunch of keywords related with the company. Combining this previously set of keywords with
412 some manually defined they have created the profiles to the companies in analyse. They used this
profiles to extract specific features from "tweets" and added to a set where there already was some
generic features. This information was further used to classify the tweets with "related/unrelated"
labels.

414 Regarding the ITC-UT [?] research group, they have, firstly, made a prediction of the company
name class according to the related-tweet ratio. After this step, a distinct heuristic was found
416 to each of the classes, using basically part-of-speech tagging and the named entity label of the
company. Their approach was a two-step classification task.

418 SINAN [?] system used an approach of ruled based heuristics, specially the existence of the
entity name both on tweets and external resources, such as Wikipedia, DBPedia ⁷ and the company
420 web-page.

The RepLab 2012 follows an identically problem as the WePS-3, the unknown-entity scenario.
422 Some research teams follow the same approach as S. Yerva et al. [?] where the use of profiles de-
scribing each company mention to correctly filter the content. DAEDALUS [?] and OXYME [?]
424 tackle a manually exploration, such as the development of dictionaries and rules sets to the detec-
tion and the classification task, and the selection of feedback terms about the entities, respectively.
426 The automatically methods were explored mainly with external resources. CIRGDISCO [?] and
ILPS [?] used the Wikipedia, while BMEDIA [?] combined it with Freebase to extract related and
428 unrelated concepts. CIRGDISCO proposed a two-step algorithm to solve the filtering task. The
first step involves the extraction of the entity related-terms from the Wikipedia and further calculus
430 of the IDF (Inverse Document Frequency) score for each term founded. The second step focus on
the idea of concept term score propagation, i.e. to propagate the labels of the high-precision clas-
432 sified tweets to the remaining, in order to increase the recall measure. ILPS tackle the filtering task
by using semanticising, where two probabilities are verified: Link Probability and Commonness.
434 The first one represents the probability that an n-gram is linked to an Wikipedia page, while the
second is the probability of an n-gram is linked to a certain concept. The ILPS group also used list
436 aggregation and disambiguation techniques to carry out this task.

At RepLab 2013, the filtering task was in a known-entity scenario where the data provided
438 to the groups consists of a collection of tweets about 61 entity names in two distinct languages,
English and Spanish. Saleiro et al. [?] have devolved POPSTAR which was the system that, using
440 a supervised learning, has obtained the best results classifying the tweets as related or non-related

⁴<http://nlp.uned.es/webs/>

⁵<http://nlp.uned.es/replab2013/>

⁶<https://wordnet.princeton.edu/>

⁷<http://wiki.dbpedia.org/>

with the entities. Their group has explored internal features (RepLab Metadata, probabilities in
442 the text, keyword similarities) and external features, such as Web Similarity (between tweet text
and the Wikipedia page text) and Freebase scores relatively to the position of the target entity in
444 the retrieved list.

The second best score in the filtering task at RepLab 2013 was obtained by V. Hangya et al. [?]
446 where their system made usage of text normalization methods, combining the textual features with
topic distribution features retrieved by a LDA (Latent Dirichlet Allocation) model. The resulting
448 features were further used in a maximum entropy classifier to perform the filtering task. The LIA
[?] group has used k-Nearest-Neighbour (kNN) algorithm with a set of discriminant features based
450 on similarity measures. They have used Bag-of-Words representation, combining TF-IDF (Term
Frequency-Inverse Document Frequency) with Gini purity criteria, for the tweets collection and
452 calculated the Jacard similarity measure.

This kind of methodologies will be a huge step to validate our dataset since it's important
454 to have only related-topic tweets to analyze the people's feelings, opinions about a correct entity
instead unrelated ones.

456 2.4.3 Topic Modeling

The emergence of topic modeling techniques was due to the people's chase of a better under-
458 standing of the available information in document corpora. Topic models provide the discovery of
certain patterns in a collection of texts and enhance specific words/terms that have a direct relation
460 to the content information [?]. There are many studies that were conducted in order to prove that
is possible to extract coherent topics from micro-blogging data using the LDA (Latent Dirichlet
462 Allocation) model [?, ?, ?]. LDA models are difficult to apply to micro-blogging texts because of
the characteristics present in this kind of text: short, mixture of contextual clues (URLs, tags, name
464 mentions with the '@'), informal language with many misspelling, acronyms and abbreviations
[?]. L. Hong et al. [?] describe Latent Dirichlet Allocation as "an unsupervised machine learning
466 technique which identifies latent topic information in large document collections". This technique
uses "bag-of-words" to each document which are represented by a probability distribution over
468 some topics, and each topic is, in turn, represented by a probability distribution over a number of
words.

470 R. Mehrotra et al. [?] have explored the improvement of the standard LDA model using
several pooling schemes of tweets, i.e. aggregating tweets by some characteristics present in its
472 content. Their polling schemes characterization range from basic scheme: where each tweet is
treated as a single document; author-wise pooling: aggregating the tweets according its author;
474 burst-score wise pooling: tweets are aggregated by the scores obtained from the execution of a
burst detection algorithm; temporal pooling: pools are formed by tweets posted at the same hour;
476 hashtag-based polling: the tweets are grouped according to its *hashtag* (#) reference, and if there
are more than one reference then the tweet is added to each of the groups. The authors evaluate the
478 resulting clusters through some metrics, such as the *purity*: verifying the average of the corrected
labeled tweets inside the clustering; *normalized mutual information* (NMI): its the calculus of the

480 matching results between the clustering and the category labels; and finally the *pointwise mutual information* (PMI): measure of the statistical independence between two words regarding the close
 482 proximity. Their approaches also were studied combining similarity tag assignment (TF and TF-IDF) and the best presented results were performed by *hashtag*-based polling with TF-similarity
 484 tag assignment regarding the purity and the NMI metrics, while the best PMI metric was obtained by the simple *hashtag*-based polling method.

486 L. Hong et al. [?] also explored the LDA models through a set of schemes formed by them.
 Their schemes diverge between user-based and term-based groups, where the user-based are ag-
 488 glomerations of messages from the same user while the term-based groups are formed by messages
 490 that have the same term in the content. In their approach they have also used Author-Topic Model
 492 which is an extension of the LDA model but the main difference is that in the LDA, each document
 494 is associated with a multinomial distribution over T topics while in the AT model the association
 is made to the author instead the document. They used JS Divergence to study the similarity be-
 tween the performed schemes. The main goal of their work was not the topic modeling but they
 proved that this sub-task can improve performances of classification, namely when the messages
 are group by the same user.

496 W. Zhao et al. [?] proposed another extension of the LDA model and named it Twitter-LDA
 498 ⁸. Their model follows the idea that each tweet is about some topic, so instead of grouping the
 tweets into schemes and than extract some topic, they tackle each tweet as a singular problem and
 extract the target of the content. In their work, the evaluation of the model was made by comparing
 500 its effectiveness against the standard LDA model and the Author-topic model. The Twitter-LDA
 results, obtained from a small set of topics in a preliminary test, have surpass the performance of
 502 the others models (standard LDA and Author-topic models).

504 The last model should be a good start to face the problem of topic modeling in our work, since
 it's open-source tool and it's available in GitHub.

2.5 Sentiment Analysis

506 Sentiment Analysis is a task of NLP (Natural Language Processing) and aims the finding of the
 polarity in opinions, sentiments of people about a specific topic contained in a document or even
 508 the overall sentiment present in it. Research done in this area has grown at an impressive pace and
 this is due to the value that this type of analysis can provide to the business world. "Marketing
 510 managers, PR firms, campaign managers, politicians, and even equity investors and online shop-
 pers are the direct beneficiaries of sentiment analysis technology" since the retrieved information
 512 can favor and make easier the decision-making process [?]. This task is composed by several dis-
 tinct problems and there are two main approaches to tackle it: supervised [?, ?] and unsupervised
 514 [?, ?]. Feldman in their work [?] enhaces the several types of problems found in the sentiment
 analysis task. One of them, the document-level sentiment analysis focus on the determination of
 516 the sentiment polarity of opinions expressed by the author in his document. Another problem that

⁸<https://github.com/minghui/Twitter-LDA>

is widely explored is the sentence-level sentiment analysis which is a deeper version of the previous. A document may have multiple opinions about a specific entity and in order to extract the polarity value about it, a phrase-level split is required. Some countermeasures must be taken into account in the polarization of phrases since the sarcasm component can be present in the content and it's very difficult to treat correctly this. There is another problem in this field named aspect-based sentiment analysis where the sentiment polarity should be directed to the aspects/topics contained in the document. In the following subsections a deep description and studied solution of this problem will be presented.

2.5.1 Lexicon-based vs Machine-learning based

Sentiment analysis in Twitter can be divided in three different approaches relatively to the sentiment classification: lexicon-based, machine-learning based or even a hybrid approach between the previous two. In the first place it's necessary to talk about what features are relevant or not in order to tackle this problem. Aggarwal et al. [?] in his work refer some of the common features used in this problem:

- **Term presence and frequency:** groups of words, named *n-grams*, and the frequency they occur in the document;
- **POS Tag:** the existence of adjectives can be relevant indicators to determine the opinion polarization;
- **Opinion words and phrases:** words that usually transmits some polarity, such as *good and bad*, or even whole phrases that don't have this type of words, e.g. "cost me an arm and a leg";
- **Negation:** the existence of negative words that may change the opinion orientation, such as "I don't like apples" which means the same as *hate*.

After the features engineering process, it may be necessary to select only a few ones to apply in the classification task. W. Medhat et al. [?] mentioned in their work some of the most used methods in this particular step:

- **Lexicon-based:** It's necessary human annotation. Starts with a small set of seed words and then a bootstrapping methodology is applied to expand the lexicon domain through the discovery of synonyms in external resources;
- **Point-wise Mutual Information:** It's a statistical method where the co-occurrence level between a given word w and a class c is computed in order to see if a feature is or not correlated with the class. The formula of calculus is given in the equation ??,

$$M_c(w) = \log\left(\frac{F(w).P_c(w)}{F(w).P_c}\right) = \log\left(\frac{P_c(w)}{P_c}\right) \quad (2.1)$$

550 where $F(w).P_c$ is the expected co-occurrence level and $F(w).p_c(w)$ is the true value of the
co-occurrence.

- 552 • **Chi-square:** It's another statistical method used to measure the correlation between the
features and the classes ??,

$$X_c^2 = \frac{n.F(w)^2.(p_c(w) - P_c)^2}{F(w).(1 - F(w)).P_c.(1 - P_c)} \quad (2.2)$$

554 where n is the total number of documents that composed the collection, $p_c(w)$ represents the
conditional probability of class c in the documents containing the word w , P_c are the fraction
of documents that contain the class c and $F(w)$ are the documents fraction that contain the
word w .

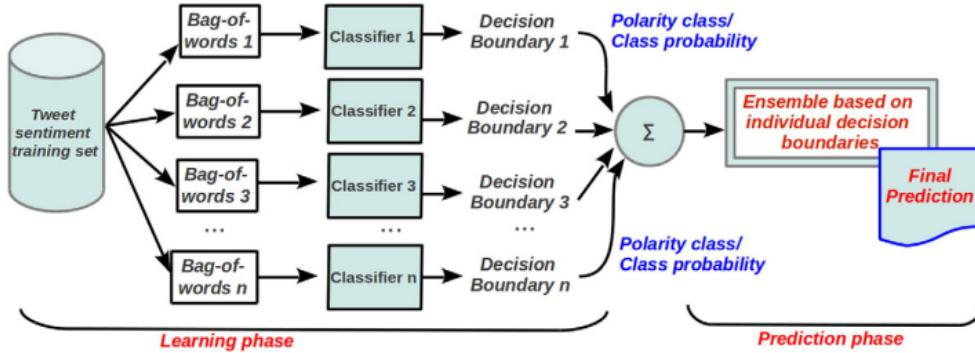
- 556 • **Latent Semantic Indexing:** It's an unsupervised method that aims the reduction of the orig-
inal set of features into a new ones through transformation techniques like PCA (Principal
Component Analysis)

560 After the conclusion of this step of features selection, the sentiment classification is conducted
and there are a very high number of techniques that can be applied.

562 A. Giachanou et al. [?] divided the **machine-learning approaches** into three different cate-
gories: supervised learning, the classifiers ensembles and deep learning.

564 Supervised learning methods focuses on the training of classification models with a man-
ually labeled dataset, also named training dataset, and various features extracted from the Twitter
566 messages in order to submit a test dataset under the model and have an automatic prediction re-
garding the polarity of the sentiment (positive, negative or neutral) in the message. There are
568 many types of classifiers, such as Naïve Bayes (NB), Maximum Entropy (ME), Support Vector
Machines (SVM), Logistic Regression (LR), Random Forest (RF) or even Conditional Random
570 Fields (CRF). In the last few years, many studies of Twitter sentiment analysis using supervised
learning were conducted. At 2009, Go et al. [?] tackle a problem of binary classification about
572 Twitter sentiment analysis using SVM, NB and ME algorithms and distant supervision to produce
their classifier models. Their dataset was composed by 1.6 million twitter messages and they don't
574 have the problem of imbalanced classes. As features they used POS tags, unigrams and bigrams
and they also have in consideration the existence of negation in the messages. The best perfor-
576 mance they have obtained was with the Naïve Bayes model and as final conclusions they said
that using POS tags as feature doesn't improve the final accuracy. Hamdan et al. [?] made some
578 experiments using SVM and NB models in Twitter messages but their set of features was different
from the previous work mentioned. The authors use DBPedia to extract concepts, WordNet to
580 extract adjectives, Senti-WordNet to extract the sentiment score of some words and also have in
consideration the existence of emoticons in the message. As final results they concluded that the
582 SVM model performance surpasses the NB model between 2%-4%, and for that they have used the
harmonic mean F-measure as evaluation metric. Saleiro et al. [?] work with Twitter messages to

Figure 2.2: The workflow of the classifiers ensembles by da Silva et al. [?]



584 study the political opinions about the five political leaders during the Portuguese bailout between
 585 2011 and 2014. The features they have used were composed by sentiment aggregate functions applying them to a non-linear regression model using the Random Forests algorithm and to a linear regression model using the Ordinary Least Squares (OLS) algorithm. The authors have grouped
 586 the dataset per month in order to see what was the monthly variation. In the validation process, a 10-fold cross validation was used and regarding to the evaluation metric they explore the Mean
 587 Absolute Error (MAE).

588 The Classifiers Ensembles approach is based on the combination of several classifiers to improve the performance of the classification task. The workflow of this approach can be verified in
 589 the figure ???. da Silva et al. [?] have explored the sentiment analysis in Twitter using a combination
 590 between Random Forests, Support Vector Machines, Multinomial Naïve Bayes and Linear Regression. The final classification decision in this kind of approach is usually made by majority voting between the models. In their work, da Silva et al. decided to calculate the average between all classification probabilities given by the models and applied that resulting value to the decision
 591 task.

592 Hassan et al. [?] have explored a different approach regarding the classifiers ensemble. They
 593 proposed a framework that was composed by seven different classification algorithms and used bootstrapping to the sample input in order to provide a small portion to each model. Their set of
 594 features had several types: semantic, POS tags, sentiment scores (SentiWordNet) and n-grams.
 595 Information Gain criteria was used to the feature selection process. By using bootstrapping in
 596 their framework, the problem of imbalance classes was reduced which could be a great advantage
 597 to explore in our approach.

598 Deep learning (DL) is the third and last method in the supervised learning approaches. DL is
 599 a recent field in the area of Machine Learning which may imply a scarcity in its addressed use to
 600 the sentiment analysis on Twitter [?]. The studies conducted in this method took advantages from
 601 the SemEval-2013 and SemEval-2014 datasets. Tang et al. [?] developed three different neural
 602 networks models to learn sentiment specific word embeddings (SSWE). The results provide by

the models were later used as features to classify the polarity of sentiment in the dataset messages.

612 The authors have combined the obtained features with others like sentiment lexicons [?], emoticons, negation, n-grams, punctuation, clusters, etc [?]. The evaluation metric used to measure
614 the performance of their classification models was the F-measure. The best result obtained to the SemEval-2013 dataset was 86.58% while for the SemEval-2014 dataset was 87.61%.

616 There are a high diversity in the methods to apply supervised learning, i.e. application of
618 machine-learning algorithms, to automatic classify the sentiment polarity either on tweets or opinion reviews. The problem in its use focuses on the features engineering which is a hard task and
620 the learning algorithms effectiveness depends on its selection. The bad choice of some features may cause that the final results obtained are not the most desirables.

Contrary to the machine-learning approaches, the **lexicon-based approaches** doesn't depend
622 on training data and features to classify the sentiment as positive, negative or neutral. In this approach the final sentiment classification is given by measuring the sentiment score of each term
624 using external resources, such as dictionaries with a large number of previously evaluated terms (SentiWordNet ⁹, SenticNet ¹⁰, LIWC ¹¹). At 2010, Thelwall et al. [?] developed a lexicon-based
626 algorithm, named SentiStrength ¹², capable of detecting the sentiment value in messages that usually have informal language, such as tweets. The algorithm have access to 298 positive terms
628 and 465 negatives and to a list of emoticons, negations and booster words to increase or decrease the sentiment value of derived words. The authors have compared their algorithm with machine-
630 learning approaches and the results were very interesting since in terms of accuracy as evaluation metric, SentiStrength has surpass the others.

632 C. Musto et al. [?] have developed a domain-agnostic framework to produce some social media analytics regarding some events that happen in Italy in the last years. They evaluate the sentiment
634 present in each tweet using lexicon-based approaches. The external resources used were SenticNet and SentiWordNet. The authors have split each tweet message by cues, such as punctuations and
636 conjunctions, creating two or more micro-phrases. After this step, each micro-phrase is classified according the scores of the terms present in the resources. The sentiment polarity of the original
638 message is obtained by summing its related micro-phrases. They also studied an emphasized approach where the Part-of-speech (POS) category of each term has a weight. Adverbs, verbs and
640 adjectives received a value greater than 1, while for the remaining categories the value was 1.

642 L. Allisio et al. [?] proposed a framework, named Felicità, in order to measure the happiness level in the Italian territory. The study was made on geotagged tweets and it was used the resources MultiWordNet ¹³ and WordNet-Affect ¹⁴. All the emoticons presented in the tweets were replaced
644 by its meaningful words. The approach used by the authors consists in for each tweet term, a search is computed in the MultiwordNet dictionary to find all the meanings the word can have.

⁹<http://sentiwordnet.isti.cnr.it/>

¹⁰<http://sentic.net/>

¹¹<http://liwc.wpengine.com/>

¹²<http://sentistrength.wlv.ac.uk/>

¹³<http://multiwordnet.fbk.eu/english/home.php>

¹⁴<http://wndomains.fbk.eu/wnaffect.html>

646 After this step, each meaning found is associated with the sentiment score present in the WordNet-Affect corpus. The sum of all meanings is calculated, assigning a value of -1, 0 or 1 to the term.
648 The tweet final classification is done by calculating the mean polarity of all terms and comparing with a heuristic constant defined by the authors.

650 The lexicon-based approaches are simpler to implement compared with the machine-learning approaches. They also presented disadvantages as the need of a continuously update of the word
652 lists (lexicon sentiment dictionaries) because the conversation themes on Twitter are always changing which may result in the absence of words in the lists, and consequently their scores [?]. For
654 this reason, the missing words are not considerate to the sentiment polarity classification and the results may not be so reliable.

656 The last approach for sentiment analysis is a mixture of the two previously presented, a **hybrid approach**. This kind of approach was explored by Ghiassi et al. [?] where they used machine-
658 learning algorithms (SVM and Dynamic Artificial Neural Networks - DAN2) with a n-gram analysis. The collection of tweets was about Justin Bieber and as features to the classifiers, the authors
660 choose emoticons, tweets that have positive and negative words, e.g. *happy or sad*, and also synonyms of this words. The model DAN2 proved be the best in the classification task. A. Kumar et
662 al. [?] mixed a log-linear regression model with lexicon-based methods. Firstly, they have made pre-processing to the tweets collection by removing the URL references, replacing emoticons with
664 their score value, calculate the percentage of caps in the message and also the sentiment orientation of the adjectives, verbs and adverbs. The overall sentiment of the tweet message was computed
666 by the linear equation of the model, which was enough to prove the efficiency of the approach explored by the author relatively to the polarity of a tweet.

668 The main advantage of the hybrid approach establishes in the no need to manually classify the dataset for its use in machine learning methods. By applying lexicon-based methods we can
670 have a labeled dataset ready to be use in ML classifiers. A disadvantage on this approach is high computational power need to bear out both approaches at the same time [?].

672 2.5.2 Aspect-based Sentiment Analysis

Aspect-based sentiment analysis (ABSA) is the most difficult problem to solve regarding the field
674 of sentiment analysis. This approach focuses on the recognition of aspects in the messages and consequently on their sentiment polarity classification.

676 In particular to the aspect extraction, an overview was already done in the subsection ?? where the topic-based approach was mentioned and described using the LDA model as the most used
678 model. There are three more approaches that we can follow to discover relevant aspects in a document: frequency-based [?], ruled-based [?] and supervised learning [?, ?].

680 The frequency-based approach focus on finding some nouns or noun phrases from a large corpus using the occurrence frequency as the main requirement. M. Hu et al. [?] used this approach
682 in order to summarize some customer reviews regarding a set of products. They used POS tagging to found nouns and noun phrases present in the document and association mining to find frequent

Background and Literature Review

684 itemsets because the features that composed the itemset are usually product features. After this
685 step, they submitted the features set to a pruning section in order to remove the meaningless ones.

686 In the ruled-based approach, S. Gindl et al. [?] have made a study in order to prove that it's
687 possible identify and extract aspects in sentences by propagating the sentiment charge to noun
688 targets following a set of defined rules. They have verified a problem when the sentiment and
689 the aspect mention are in different sentences. Simple propagation rules are not enough to found
690 the aspects. To overcome this, they have defined another rule where if a sentence starts with a
691 pronoun, the target aspect is probably the last noun identify in the previous sentence.

692 Supervised learning approaches are based in sequential learning methods, such as Conditional
693 Random Fields (CRF) and Hidden Markov Models (HMM). The mentioned methods are similar
694 because both of them attempt to discover patterns relative to an input data set. These types of
695 methods are often used in the aspect extraction task of opinion mining. N. Jakob et al. [?] has built
696 a classification model using the CRF algorithm in order to enhance the target of each opinion in the
697 reviews. The domain of their dataset was constituted by four independently categories: movies,
698 web-services, cars and cameras. Since the approach taken by the authors was through machine
699 learning classifiers, they had the need of establishing a set of features to train the model. The
700 features used for the classification vary from the string of each token, the POS tag for each token,
701 the level of dependency between each token and the opinion expressed as well as its word distance,
702 and, finally, the last feature is the opinion sentence itself to allow the CRF algorithm the ability of
703 distinguishing when a token is present or not in a sentence that is an opinion. Regarding the system
704 validation, the authors applied also a 10-fold cross-validation to see if the model performance
705 improves or not. As performance metrics to evaluate the system, they used the precision (Equation
706 ??),

$$\frac{TP}{TP + FP} \quad (2.3)$$

the recall (Equation ??)

$$\frac{TP}{TP + FN} \quad (2.4)$$

708 and the F-measure (Equation ??) that is the harmonic mean between the previous two metrics.

$$2. \frac{precision \cdot recall}{precision + recall} \quad (2.5)$$

W. Jin et al. [?] also worked under the opinion reviews and tried to find relevant opinion targets
710 in the content. They have developed a novel framework based on machine learning techniques.
711 The framework, named OpinionMiner, appeals to a classification model with the HHM algorithm
712 and to a bootstrapping approach in the training step. The bootstrapping divides the main process
713 of training in two sub-process as well as the training dataset into little portions in a randomly way.
714 Each sub-process has his own HHM model and after its training step, the main process only selects

the objects which their label is agreed by both classifiers. The bootstrapping process is repeated
716 until no more targets in the objects can be discovered.

Regarding the polarity sentiment classification, the majority of the works studied appeals to
718 one of the approaches described in the section ??.

2.6 Conclusion

720 This chapter had the objective of review some basic concepts that may be relevant to contextualize
the reader about the problem of performing analysis in social media streams, e.g. Twitter mes-
722 sages. Hence, the literature studied was divided in several points in order to have a overview about
what is already done and what are the approaches that some author proposed to tackle each of the
724 sub-problems that composed the main problem in this dissertation work. After a careful research,
it was possible to identify that there are a great diversity of approaches to each sub-problem,
726 whether it be disambiguation, filtering, topic detection or even sentiment analysis.

An important point identified in the literature was the few works done using deep learning to
728 take the problem of sentiment analysis with supervised leaning approaches, since its applicability
in the artificial intelligent field has grown at an exponential level in the recent years.

730 Regardless the task that the authors dealt with, it was possible to identify that the features en-
gineering process and its selection, when their proposed solution used classifier models, is similar.
732 This may be an advantage to the development of the different modules that composed the proposed
framework in this dissertation work.

734 The framework modules will also have classifier models, so the evaluation and validation of it
is important. The literature review shows a large set of evaluation metrics to do this step.

736 In short, it is expected from the reader that this review to the State-of-the-Art has provided
a coherent understanding regarding the study of different real scenarios using the social media
738 streams as source of information.

Chapter 3

Framework

⁷⁴⁰ 3.1 Requirements

3.2 Architecture Overview

⁷⁴⁴ 3.3 Data Collection

3.4 Data Pre-processing

⁷⁴⁶ 3.5 Text Analytics

3.5.1 Word Embeddings

⁷⁴⁸ Mikolov et al. [?] has developed a powerful computational method named *word2vec*. The method
is capable of learning distributed representations of words, each word being represented by a dis-
⁷⁵⁰ tribution of weights across a fixed number of dimensions. Authors have also proved [?] that this
kind of representation is robust when encoding syntactic and semantic similarities in the embed-
⁷⁵² ding space.

The training objective of the skip-gram model, as defined by Mikolov et al. [?], is to learn the
⁷⁵⁴ target word representation, maximizing the prediction of its surrounding words given a predefined
context window. For instance, to the word w_t , present in a vocabulary, the objective is to maximize
⁷⁵⁶ the average log probability:

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log P(w_{t+j} | w_t) \quad (3.1)$$

where c is the size of the context window, T is the total number of words in the vocabulary and
⁷⁵⁸ w_{t+j} is a word in the context window of w_t . After training, a low dimensionality embedding matrix

E encapsulates information about each word in the vocabulary and its use (i.e. the surrounding contexts).
760

Later on, Q. Le and Mikolov [?] developed paragraph2vec, an unsupervised learning algorithm
762 operating on pieces of text not necessarily of the same length. The model is similar to word2vec
but learns distributed representations of sentences, paragraphs or even whole documents instead of
764 words. We used *paragraph2vec* to learn the vector representations of each tweet and tried several
configurations in the model hyperparameters.

766 **3.5.2 Latent Dirichlet Allocation (LDA)**

D. Blei et al. [?] have developed a generative statistical model that makes possible the discover
768 of unknown groups and its similarities over any dataset. The model tries to identify what topics
are present in a document by observing all the words that composing it, producing as final result
770 a topic distribution. An interesting point in this model is that the only features it analyses are
the words passing through the training process. The model takes into consideration two different
772 distributions - distribution of words over topics and distribution of topics over the documents -
being each document seen as a random mixture of topics and each topic as a distribution of words.

774 **3.6 Visualization**

Chapter 4

⁷⁷⁶ Experiments

The developed framework presented and described in Chapter ?? oblige to the validation of each
⁷⁷⁸ module in order to assure consistency and robustness in the results that such system produced.
Having this considered, we stipulate several experiments, in which each of them was related to a
⁷⁸⁰ specific task.

⁷⁸² 4.1 Exploratory Data Analysis

The main goal of this section is the devise of relevant analysis taking into consideration the five
⁷⁸⁴ different collected datasets. Since this dissertation is supported in experiments using real-world
data, such analysis is crucial in order to gain better knowledge of the intrinsic characteristics of
⁷⁸⁶ it. A tweet provides some fields of interest, such as, the text message, date of creation, language,
and the *entities*, which are constantly analysed in several data analytics systems. An *entity* is
⁷⁸⁸ metadata and additional contextual information contained in the tweet and is composed by the
hashtags, *user mentions*, *urls* and *media* fields. We count the amount of tweets containing this
⁷⁹⁰ kind of information for all the cities, London, New York, Melbourne, Rio de Janeiro and São
Paulo, and projected some data visualizations for different temporal frequencies. The following
⁷⁹² subsections are divided into three different categories: (1) Geographical Distribution, (2) Temporal
Frequencies and (3) Metadata Composition. Additionally, we discuss the results of each city, as
⁷⁹⁴ well as the main observable differences.

4.1.1 Geographic Distributions

⁷⁹⁶ The bounding-boxes selected to perform the data collection had its origin in an open-sourced tool

Experiments



Figure 4.1: Search Bounding-boxes for the data collection

4.1.2 Temporal Frequencies

4.1.3 Metadata Composition

Table ?? present the results for the *entities* count and it is possible verify that *urls* and *user mentions* are most used ones over both datasets. Such information shows that users tend to tag another ones in a message meaning that tweets are used as a mean of communication.

Table 4.1: Datasets entities statistics

City	Hashtags (#)		User Mentions (@)		URLs		Media	
	Total	%	Total	%	Total	%	Total	%
Rio de Janeiro	525,550	5%	1,340,334	13%	1,509,742	14%	389,864	4%
São Paulo	585,365	12%	1,072,566	22%	885,369	18%	302,579	6%

The grouping of tweets by the day of the week is illustrated in Figure ?? and particular points can be observed and considered uncommon. For both cities distributions and with respect to geo-located tweets, Friday is the less active day while the major activity occurs in the first days of the week. This is a strange phenomenon since Friday is transition of the labour week days to the weekend and people could use more the microblog service to share their free time.

4.2 Topic Modelling

This section is related to the experiment of automatically characterize tweets in two different Brazilian cities, Rio de Janeiro and São Paulo. We used an unsupervised learning approach to

Experiments

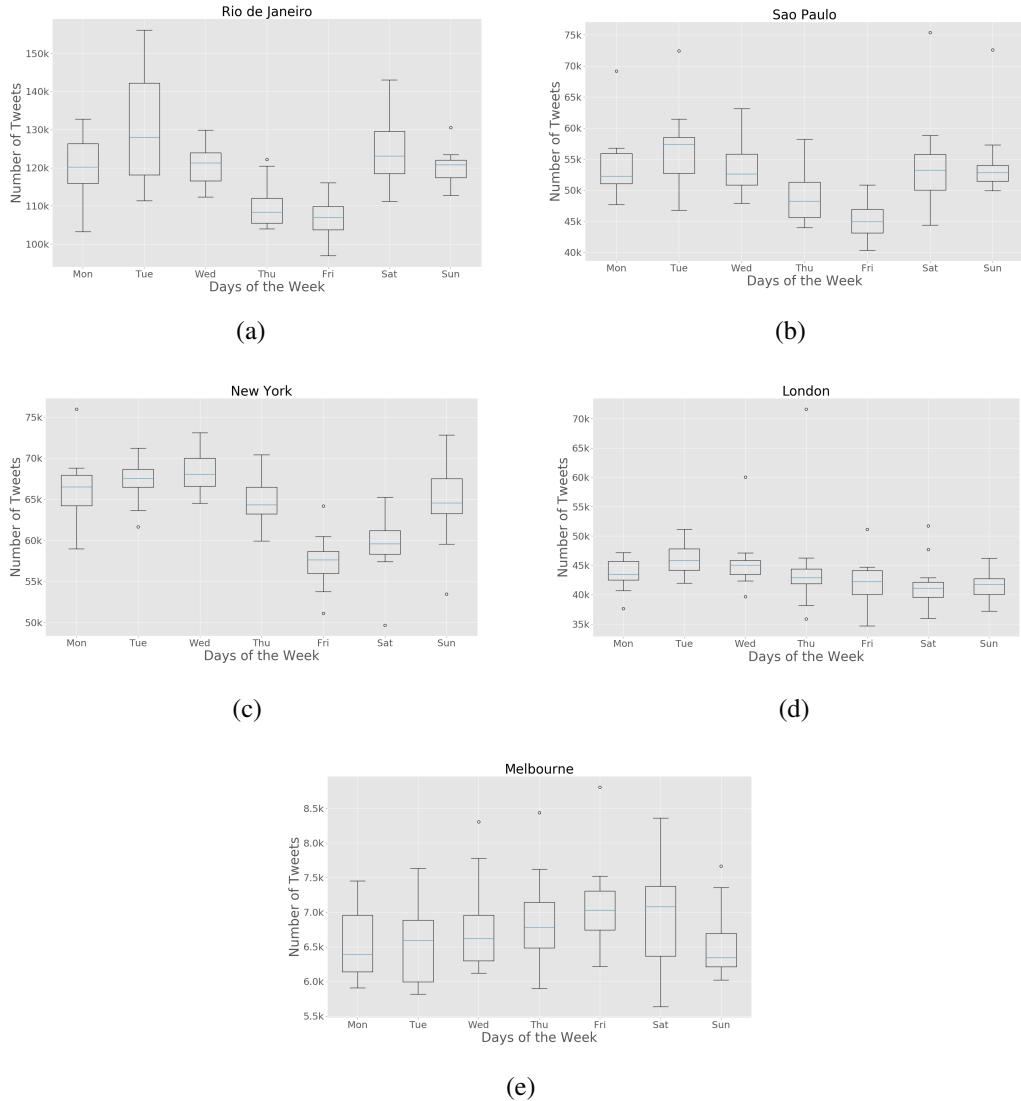


Figure 4.2: Days-of-the-week box-plots for the volume of tweets (a) Rio de Janeiro (b) São Paulo (c) New York City (d) London (e) Melbourne

tackle the task of topic modelling in order to compare both cities and see if there are differences between subjects people talked about. Automatic characterization of text messages is a laborious and time consuming task since it is necessary to assure the right level of abstraction in the learning model; very much similarly to human minds, which essentially present a bounded rationality nature, our learning model needs to be trained in order to assimilate the necessary knowledge and perform the appropriate analogies so as to discover different topics within the tweets' contents. The premises to implement such a mechanism are presented and discussed in the following subsections.

Experiments

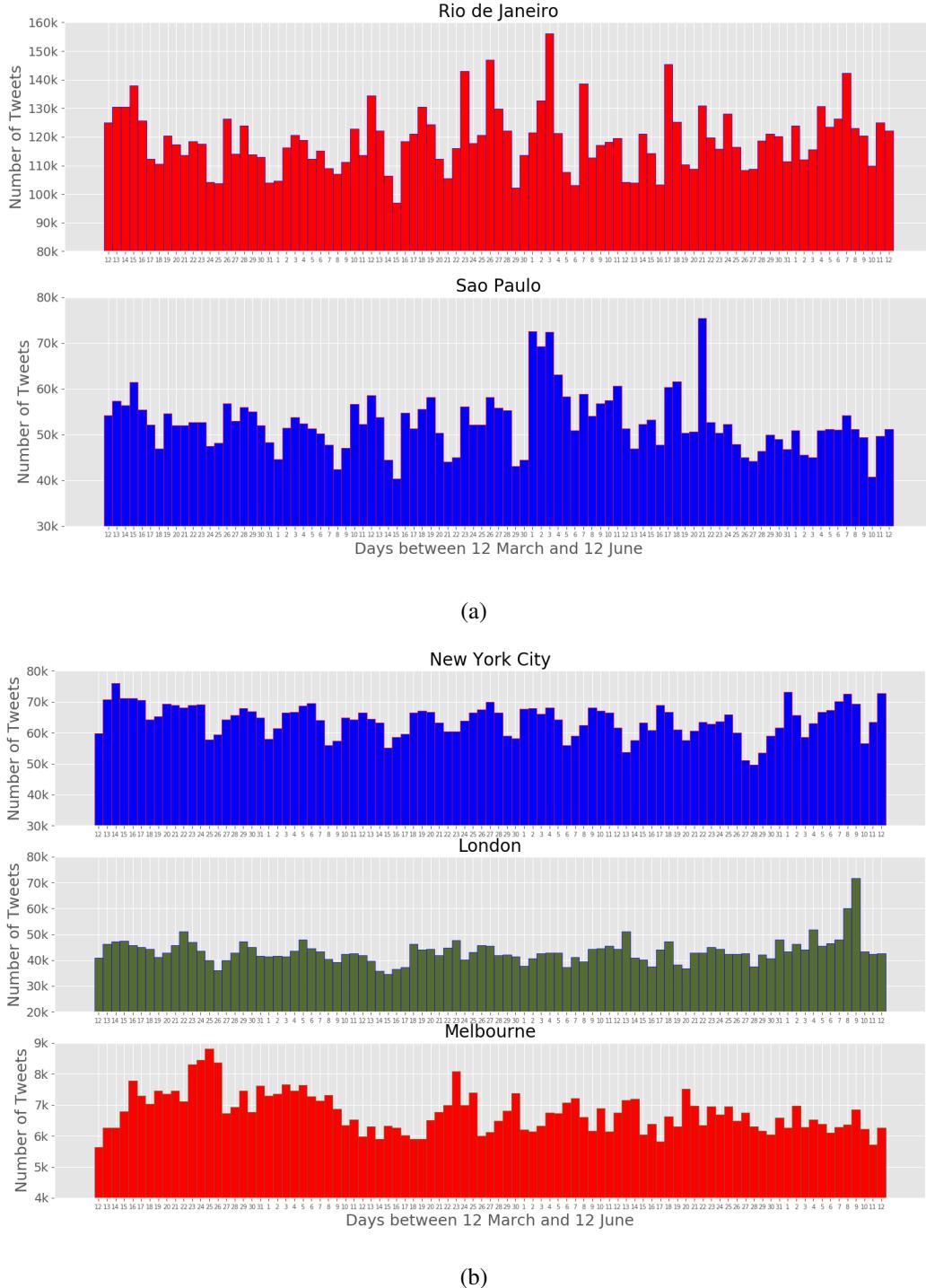


Figure 4.3: (a) Rio de Janeiro and São Paulo - Portuguese Cities (b) New York City, London and Melbourne - English Cities

818 **4.2.1 Data Selection**

The data selected to conduct this experiment is correspondent to a period of two months, between
820 days March 12 and May 12, 2017.

The resulting datasets sum up a total of 12.5M and 6.3M tweets for Rio de Janeiro and for São
822 Paulo, respectively. Due to the problem detected in Section ??, we filtered the data in order to only
use the tweets that were actually inside the cities' areas. The final composition of the datasets is
824 presented in Table ??, and the results of the filtering process shown that almost 6M tweets were
not located inside the bounding-boxes of the cities.

Table 4.2: Datasets composition

City	All	PT	Non-PT	In Bounding-Box	Out Bounding-Box	PT and In Bounding-Box
Rio de Janeiro	12,531,000	10,570,000	1,961,000	8,644,000	3,886,000	7,353,000
São Paulo	6,352,000	4,886,000	1,466,000	4,247,000	2,105,000	3,313,000

The subset of data composed by Portuguese tweets and located inside the cities' bounding-boxes was used to conduct the experiment described in this section. Such subset can be sum up to
826 a total of 7.3M and 3.3M for Rio de Janeiro and São Paulo, respectively.
828

4.2.2 Data Preparation

Usually, to tackle topic modelling tasks in text documents it is required several pre-processing
830 steps. Such pre-processing to the data helps the operations made by the LDA model, which is the
technique used here. Removing unnecessary words, transforming words into their root form as so
832 deleting all the punctuation are some of the common text mining pre-processing steps. Here, each
834 tweet of both datasets was submitted to a required group of pre-processing operations in order to
train a LDA model and proceed with the experiments. The pre-processing steps were the ones
836 detailed below.

- **Lowercasing:** Every message presented in a tweet was converted into lower case;
- **Cleaning Entities and Numbers:** Removing *URLs*, user mentions, *hashtags* and digits from the text message;
- **Lemmatization:** Only plural words were transformed into singular ones;
- **Transforming repeated characters:** Sequences of characters repeated more than three times were transformed, e.g. "loooool" was converted to "loool";
- **Punctuation Removal:** Every punctuation was removed as well as smiles (e.g. :), :-), =D) or even *emoticons*;
- **Stop Words Removal:** The removing of this kind of words was made using the Portuguese NLTK dictionary;

- **Short Tokens Removal:** Words such as 'kkk', 'aaa', 'aff' and other of the same style were removed.

After the data preparation phase, 772,017 tweets have their message empty which conclude that its content was irrelevant for the final experiment phase.

4.2.3 Features Selection

Topic modelling requires, like in other learning model, a group of features to be trained. In this case, we used the Bag-of-Words representation matrix - which is a representation where each document is converted to a frequency vector according to the number of occurrences of each word in the message. The set of features was limit to a dictionary containing 10,000 words and it only took into account uni-grams in the message content. The dictionary was also limited to words that occur in a maximum percentage of 40% in the whole dataset, avoiding common words that were not removed because they were not included in the NLTK Stop Words list. The minimal occurrence value for a word being considered was set to 10.

4.2.4 LDA Model Parametrization

In order to understand and see the LDA model performance, we set five different numbers for the topics results parameter of the training process: 5, 10, 20, 25 and 50 topics, being this the one with better results. The number of iterations to train the model was set to 20, since our desired was to reproduce the experiment made by G. Lansley et al. [?] to the city of London. Finally but not the least, each tweet in the datasets was treated as a single document comprehending that, in total, 6,580,983 different documents were used in the model training process. The complete pipeline according to all the steps taken to conduct this experiment is observable in Figure ??.

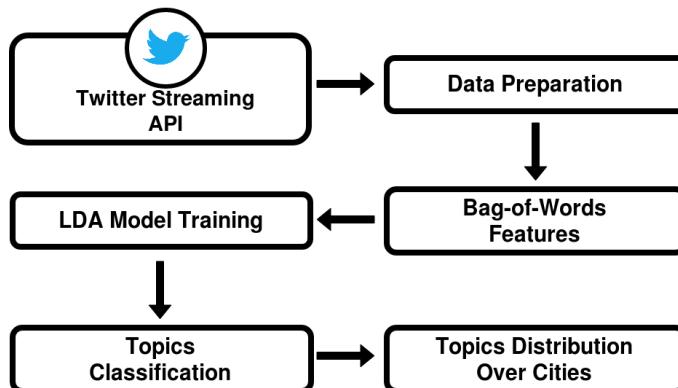


Figure 4.4: Correspondent pipeline of the topic modelling experiment

4.2.5 Results and Analysis

To evaluate the experimental results obtained for each model (where the difference underlies on the variation of the number of topics), a list with the most frequent 50 words for each topic was

Experiments

extracted. In Table ?? we can observe a sample (20 top words) selected out of the 50 studied.
 872 Nonetheless, the final evaluation takes into consideration the 50 frequent words.

Table 4.3: Example of the topics classification

Words (only 20 words)	Topic Classification
paulo, vai, hoje, dia, jogo, ser, melhor, time, vamo, brazil, todo, santo, brasil, gol, cara, aqui, agora, corinthiam, ano, palmeiro, vem, ...	Sports and Games
vou, dia, dormir, queria, hoje, ficar, casa, semano, quero, ter, ainda, hora, agora, sono, aula, acordar, acordei, cedo, fazer, prova, ...	Wake-up Messages
top, social, artist, vote, the, award, army, bom, voting, doi, bogo, oitenta, sipda, today, vinte, prepara, cypher, oito, quatro, man, ...	Voting and Numbers
marco, nada, falar, emilly, gente, quer, nao, pessoa, nunca, fala, vai, falando, sobre, chama, agora, manda, vem, mensagem, vivian, bbb, ...	Big Brother Brazil 2017
paulo, brazil, sao, santo, vila, just, parque, posted, photo, shopping, paulista, centro, bernardo, jardim, cidade, avenida, praia, santa, campo, academia	Tourism and Places

We also selected and manually analyse a random sample (with the size of 200) of tweets for
 874 each topic. This sampling was done in order to get better consistency and trustiness about the
 classification and characterization of the tweets.

876 It was found a group of 50 topics which had the largest number of distinct topics between
 them. However, there were topics which theme was the same (e.g. Love and Romance Problems
 878 or Brazilian Football *versus* European Football). Within this, such groups were aggregate into the
 same topic, *Relationships* and *Sports and Games*, respectively. After this grouping process, a total
 880 of 29 different topics was achieved.

Some tweets that have added complexity to our classification objective, such as, for example,
 882 "queria namorar um mano parecido com o josh" (Relationship) and "como eu queria meus amigos
 aquia agora cmg" (Friendship), raised some doubts about which topic this tweets may belong:
 884 Relationship, Friendship or even Actions or Intentions. In a perspective of context, the first tweet
 belongs to the theme *flirt*, which is directly related to Relationship. The theme on the second
 886 tweet is missing the company of friends, i.e. conviviality, which is related to Friendship. The
 decision of join the two topics was due to the proximity between them which have as content both
 888 types of tweets, talking about love/relationship and friendship, and with this in consideration both
 topics should be aggregated in order to assure the desired consistency in the classification.

890 The final set of topics (50 topics) to be considered was selected accordantly to the most re-
 curring subjects. The final classification and details associated with the whole dataset for each
 892 city is presented in Table ???. Almost every topics demonstrated a balanced distribution, with ex-
 ception of *Relationships and Friendship* and *Personal Feelings* for Rio de Janeiro and São Paulo,
 894 respectively. The difference that appear in this topics is a consequence of the final grouping pro-
 cess, since there was a considerable number of words been shared among this topics. This issue
 896 complicated our classification task, compelling to an high amount of undesired aggregations.

Additionally to the manual verification of a sample of tweets for each topic, we also produced
 898 a temporal week day distribution, with the objective to observe if some topics had more mentions

Experiments

Table 4.4: Final results of the LDA topics aggregation

Topic Group	Rio de Janeiro		São Paulo		Diff (%)
	No. Tweets	Percentage (%)	No. Tweets	Percentage (%)	
Academic Activities	101,590	1,54%	90,616	3,30%	-1,76%
Actions or Intentions	600,030	9,12%	128,710	4,69%	+4,43%
Anticipation and Socialising	132,606	2,01%	0	0,00%	+2,01%
BBB17	122,054	1,85%	68,385	2,49%	-0,64%
Body, Appearances and Clothes	160,342	2,44%	71,447	2,60%	-0,17%
Food and Drink	167,204	2,54%	58,407	2,13%	+0,41%
Health	119,013	1,81%	0	0,00%	+1,81%
Holidays and Weekends	104,695	1,59%	79,610	2,90%	-1,31%
Informal Conversations	272,502	4,14%	138,848	5,06%	-0,92%
Live Shows, Social Events and Nightlife	359,342	5,46%	140,240	5,11%	+0,35%
Mood	139,287	2,12%	138,399	5,04%	-2,92%
Movies and TV	285,198	4,33%	39,778	1,45%	+2,89%
Music and Artists	84,407	1,28%	78,142	2,85%	1,56%
Negativism, Pessimism and Anger	229,104	3,48%	183,050	6,67%	-3,18%
Numbers, Quantities and Classification	86,897	1,32%	78,160	2,85%	-1,53%
Optimism and Positivism	106,714	1,62%	39,725	1,45%	+0,18%
Personal Feelings	375,735	5,71%	532,331	19,38%	-13,67%
Politics	81,254	1,23%	46,758	1,70%	0,47%
Relationships and Friendship	1,524,804	23,17%	187,541	6,83%	+16,34%
Religion	183,174	2,78%	66,788	2,43%	+0,35%
Routine Activities	334,216	5,08%	82,421	3,00%	+2,08%
Slang and Profanities	241,676	3,67%	44,620	1,62%	+2,05%
Social Media Applications	105,809	1,61%	44,073	1,60%	+0,01%
Sport and Games	382,479	5,81%	133,047	4,84%	+0,97%
Tourism and Places	59,288	0,90%	86,519	3,15%	-2,25%
Transportation and Travel	130,261	1,98%	63,923	2,33%	-0,35%
Weather	91,302	1,39%	42,588	1,55%	-0,16%
Shopping	0	0,00%	44,470	1,62%	-1,62%
Voting	0	0,00%	37,687	1,37%	-1,37%

in certain days than others.

900 For making such observations some assumptions were made in relation with some *hot* topics.
 More specifically, we think that is valid to assume that people will talk more about *Religion* in the
 902 weekend, since they go to the church in those days. The same result is likely to happen for topics
 like *Holidays and Weekends* or *Sports and Games*, since events related to this thematics occur
 904 during specific time-frames.

Only 12 topics of the finals 29 were selected for this part of the study, predicting them and
 906 comparing the final results, such as, but not limited to, *Sports and Games*, *Religion*, *Holidays and*
Weekends, *Movies and TV*, *Live Shows*, *Social Events and Nightlife*. The temporal distribution is
 908 showed in Figure ?? as a heat map, where each row is independent from the others.

The necessity of applying such restrictions is due to the need of seeing in which days each topic
 910 is more talked about. For both cities the topic *Sports and Games* is more mentioned in Tuesdays
 and Saturdays. Indeed, this observation correlates with the days that topic-related events happens.
 912 Namely, Tuesdays and Wednesday correspond to the days when the *UEFA Champions League*
 competition happens and Saturdays and Sundays to the days of *Brazilian Football League* games.
 914 *Holidays and Weekends* was a topic with interesting results regarding the temporal distribution,
 presenting Sundays as the day where more people talk about it.

916 Furthermore, it is worth mentioning that our model had successfully discover a topic related to

Experiments

Big Brother Brazil 2017 (BBB17), a well-known reality show. The amount of geo-located tweets concerning this topic was considerable (1.85% and 2.49%, in RJ and SP, respectively), rising the question about what led people to geo-located them in such topic.

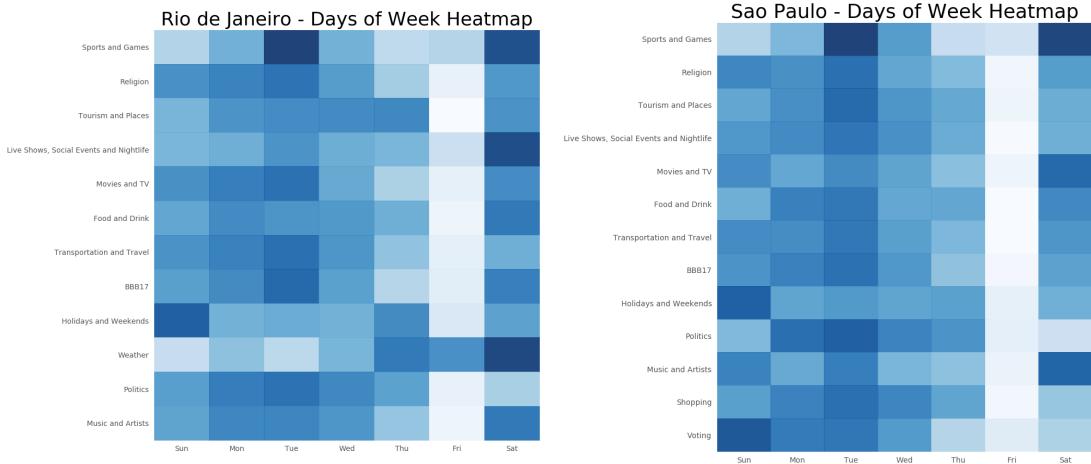


Figure 4.5: Day-of-the-week activity per each topic in both cities, Rio de Janeiro and São Paulo

4.2.6 Final Remarks

The methodology reported across this experiment is concerned with topic modelling over two datasets from two Brazilian cities in order to characterize the topics that people talked about and compare the results in both scenarios. LDA models usually require documents of large size, or at least more complex than a single tweet, in order to get good performance. A traditional approach was followed considering each tweet as a document instead of trying aggregate tweets in more complex documents taking into consideration some criteria, e.g. grouping by date and hour. The final results showed that topics in both cities are very similar and only two of them are unique. With exception of topics - *Relationships and Friendship* and *Personal Feelings*, the percentage difference between similar topics was comprehended in the interval 0.16-4.43% evidencing the fact that both cities are similar besides the different factors that characterize each one: population, culture, lifestyle and also the region where the city is located in. Although all this analysis, we can not assure that inside a topic we do not have more topics hidden. Our classification was limited to the verification of the 50 top words and the manually verification of a sample of 200 tweets since the resulting amount of tweets for each topic is impossible to verify one by one. Due to this, another classification approach need to be explored and a promising one was proposed by D. Ramage et al. [?]. The classification will be automatic by adding a supervised extra layer to the pipeline. However, to assure trustiness in the results the data may be manually labelled for the training phase of the model classification or, at least, have reliable sources, for example, exploring the topics provide by the Wikipedia articles¹.

¹<https://dumps.wikimedia.org/ptwiki/20170601/>

940 4.3 Travel-related Classification

The main goal of this section is to detail the experiment that supports the characterization of travel-related tweets in Rio de Janeiro and in São Paulo. Considering the volume of the collected data, it was then necessary to automatically identify tweets whose content somehow suggests to be related to the transportation domain. Conventional approaches would require us to specify travel-related keywords to classify such tweets. On the contrary, our approach consisted in training a classifier model to automatically discriminate travel-related tweets from non-related ones.

One big challenge always present in text analysis is the sparse nature of data, which is especially the case in Twitter messages. Conventional techniques such as Bag-of-Words tend to produce sparse representations, which become even worse when data is composed by informal and noisy content.

Word embeddings, on the other hand, is a text representation technique that tries to capture syntactic and semantic relations from words. The result is a more cohesive representation where similar words are represented by similar vectors. For instance, "*taxi*"/"*uber*", "*bus/busão/ônibus*", "*go to work*"/"*go to school*" would yield similar vectors respectively. We are particularly interested in exploring the characteristics of word embeddings techniques to understand which extent it is possible to improve the performance of our classifier to capture such travel-related expressions. In the following subsections, we describe the necessary steps to build our classification model.

958 4.3.1 Data Preparation

Each tweet of our training and test sets was submitted to a small and basic group of pre-processing operations, as detailed below.

- **Lowercasing:** Every message presented in a tweet was converted into lower case;
- **Transforming repeated characters:** Sequences of characters repeated more than three times were transformed, e.g. "loooool" was converted to "loool";
- **Cleaning:** URLs and user mentions were removed from the text.

4.3.2 Features Selection

966 We established the use of different groups of features to train our classification model, namely bag-of-words, bag-of-embeddings - word embeddings dependent technique - and both combined.
 968 Such groups are detailed below.

- **Bag-of-words (BoW):** This group of features was obtained using unigrams with standard bag-of-words techniques. We considered the 3,000 most frequent terms across the training set excluding the ones found in more than 60% of the documents (tweets);

- 972 • **Bag-of-embeddings (BoE):** We applied bag-of-embeddings to each tweet using a *doc2vec*
 974 model ² combining Deep Learning and *paragraph2vec*. The model was trained with 10
 976 iterations over the whole Portuguese dataset using a context window of value 2 and feature
 978 vectors of 100 dimensions. We then took the corresponding embedding matrix to yield the
 group of features fed into our classification routine.
- **Bag-of-words plus Bag-of-embeddings:** We horizontally combined both the above matri-
 ces into a single one and used it as a single group of features.

4.3.3 Training and Test Datasets

- 980 The construction of the training and test sets followed a traditional approach. We thus tried to
 982 select balanced training sets, to which it was necessary to identify tweets that could possibly be
 travel-related. We were inspired by a strategy used in the study by Maghrebi et al. [?], which
 984 consists in searching tweets from a collection using specific travel terms and regular expressions.
 986 Using the terms declared in Table ?? combined with the regular expression *space + term + space*,
 we found about 30,000 tweets. From this subset, we randomly selected a small sample of 3,000
 988 tweets to manually confirm if they were indeed related to travel topics. After this manual annota-
 tion we selected 2,000 tweets and used them as positive samples in the training dataset.
- 988 In order to select negative samples for the training dataset we randomly selected 2,000 tweets
 and also manually verified their content to assure that they were not travel-related. Finally, our
 990 training set was composed by 4,000 tweets, from which 2,000 were travel-related and 2,000 were
 not. We selected 1,000 tweets randomly that were not present in the training set to build the test
 992 set, and then manually classified them as travel-related or non-travel-related. In the end, 71 tweets
 were found to be travel-related and whereas 929 were not.

Table 4.5: Travel terms used to build the training set

Mode of Transport	Terms
Bike	bicicleta, moto
Bus	onibus, ônibus
Car	carro
Taxi	taxi, táxi
Train	metro, metrô, trem
Walk	caminhar

994 4.3.4 Estimators and Evaluation Metrics

Support Vector Machines (SVM), Logistic Regression (LR) and Random Forests (RF) were the
 996 classifiers used in our experiments. The SVM classifier was tested under three different kernels,
 namely *rbf*, *sigmoid* and *linear*; the latter proved to obtain the best results.

²<https://radimrehurek.com/gensim/models/doc2vec.html> (last visited on 9 June, 2017)

Experiments

998 The LR classifier was used with the standard parameters, whereas the RF classifier used 100 trees in the forest. The gini criterion and the maximum number of features were limited to those
1000 as aforementioned in Section ??, in the case of the RF classifier.

1002 To evaluate the performance of the classifiers in our experiences we used five different metrics.
1002 Firstly we compute a group of three per-class metrics, namely precision, recall and the F1-score.
1004 Bearing in mind this study considers a binary classification, metrics were associated with the travel-related class only, i.e. the positive class. Therefore, the interpretation for each metric is provided below:

- 1006 • **Precision:** Represents the fraction of correct predictions for the travel-related class (Equation ??).
- 1008 • **Recall:** Represents the fraction of travel-related tweets correctly predicted (Equation ??).

$$Precision = \frac{tp}{tp + fp} \quad (4.1) \qquad \qquad \qquad Recall = \frac{tp}{tp + fn} \quad (4.2)$$

1010 where tp is related to the true positives classified tweets, fp represents the false positives and fn are the false negatives.

- **F1-score:** Represents the harmonic mean of precision and recall.

$$F1_{score} = 2 * \frac{precision * recall}{precision + recall} \quad (4.3)$$

1012 Once these first three metrics only showed us the performance of the classifier for a discrimination threshold of 0.5, we decided to calculate another metric. The ROC (Receiver operating
1014 characteristic) curve gives us the TPR (True positive rate) and the FPR (False positive rate) for all possible variations of the discrimination threshold. Through the ROC curve, we compute the
1016 area under the curve (AUC) to see what was the probability of the classifier to rank a random travel-related tweet higher than a random non-related one.

1018 4.3.5 Results and Analysis

Table ?? presents the results obtained using the different features combination for our test set composed by 1,000 tweets manually annotated. According to the evaluation metrics we conclude that the bag-of-word and bag-of-embeddings combined produced better classification models. The model produced by the Linear SVM performed slightly better than the LR and the RF. Interesting to note is that BoW features have influence on the precision scores obtained from our results, producing more conservative classifiers. Regarding the recall results, we can see that the Logistic Regression using only bag-of-embeddings features was the model with best results; perhaps if the precision is taken into consideration, the same conclusions will not be possible. Analysing the

Experiments

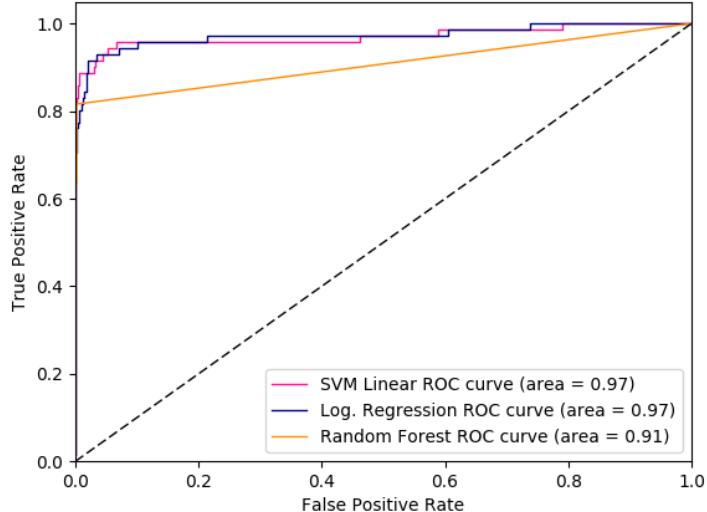
Table 4.6: Classifiers Experiences

Classifier	Features	Precision	Recall	F1-score
Linear SVM	BoW	1.0	0.6761	0.8067
	BoE	0.4338	0.8309	0.5700
	BoW + BoE	1.0	0.7465	0.8548
Logistic Regression	BoW	1.0	0.6338	0.7759
	BoE	0.4444	0.8451	0.5825
	BoW + BoE	1.0	0.6761	0.8067
Random Forest	BoW	1.0	0.6338	0.7759
	BoE	0.2298	0.8028	0.3574
	BoW + BoE	1.0	0.6338	0.7759

scores provided in Table ??, the best model under the F1-score was the Linear SVM, with a score
1028 of 0.85.

The performance of all three classifiers is illustrated using the ROC Curve in Fig. ??.
1030 The area under the curve of the Receiver Operating Characteristic (AUROC) was very similar for both the
1031 Logistic Regression and the Linear SVM models. The results obtained from the Random Forest
1032 model were not so promising as expected.

Figure 4.6: ROC Curve of SVM, LR and RF experiences



After the selection of our classification model, we decided to classify all the Portuguese dataset
1034 and draw some statistics from the results. The trained Linear SVM classifier was used to predict
1035 whether tweets were travel-related or not, since it was the model presenting the best score under
1036 the F1-score metric (as shown in Table ??). From a total of 7.8M tweets, our classifier was able
1037 identified 37,300 travel-related entries.

Fig. ?? depicts the distribution of travel-related tweets over the days of the week. We can see
1038 that the first three business days (Monday, Tuesday and Wednesday) are the ones on which the
1040 Twitter activity is higher for both cities in our study.

Experiments

Figure 4.7: Positive Predicted Tweets per Day of Week

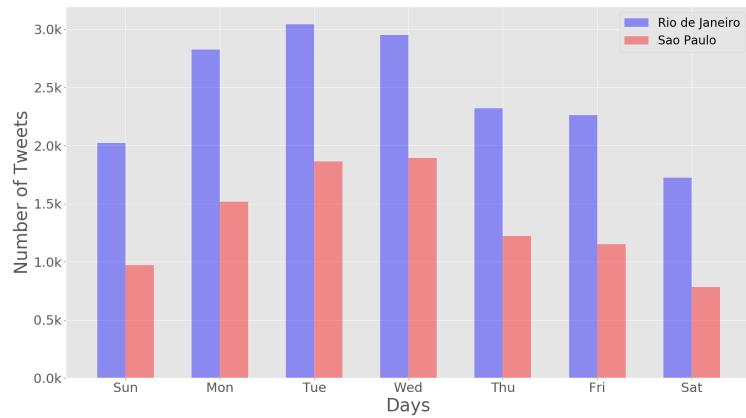
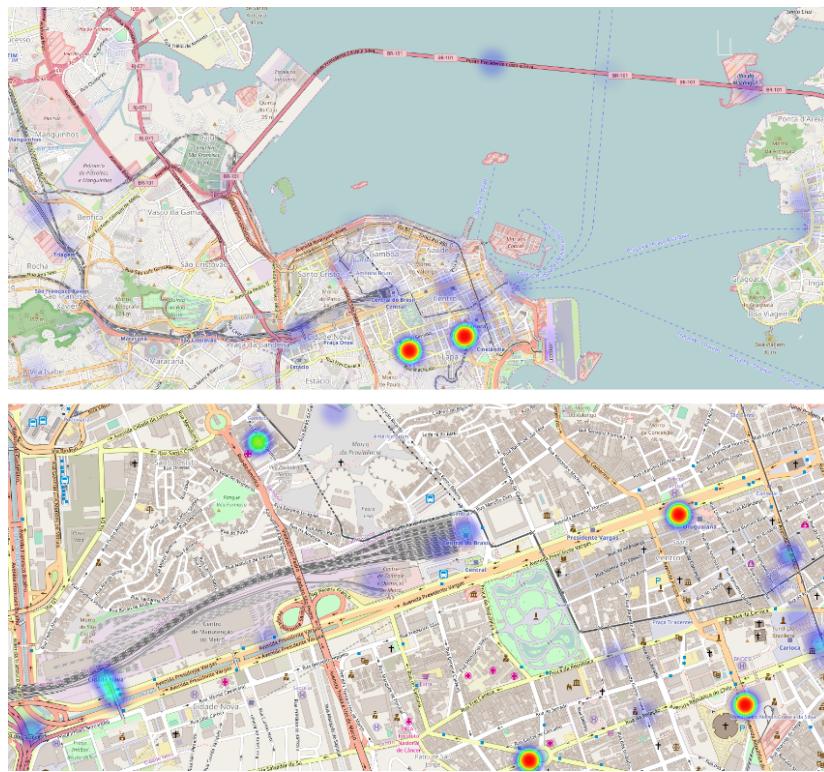
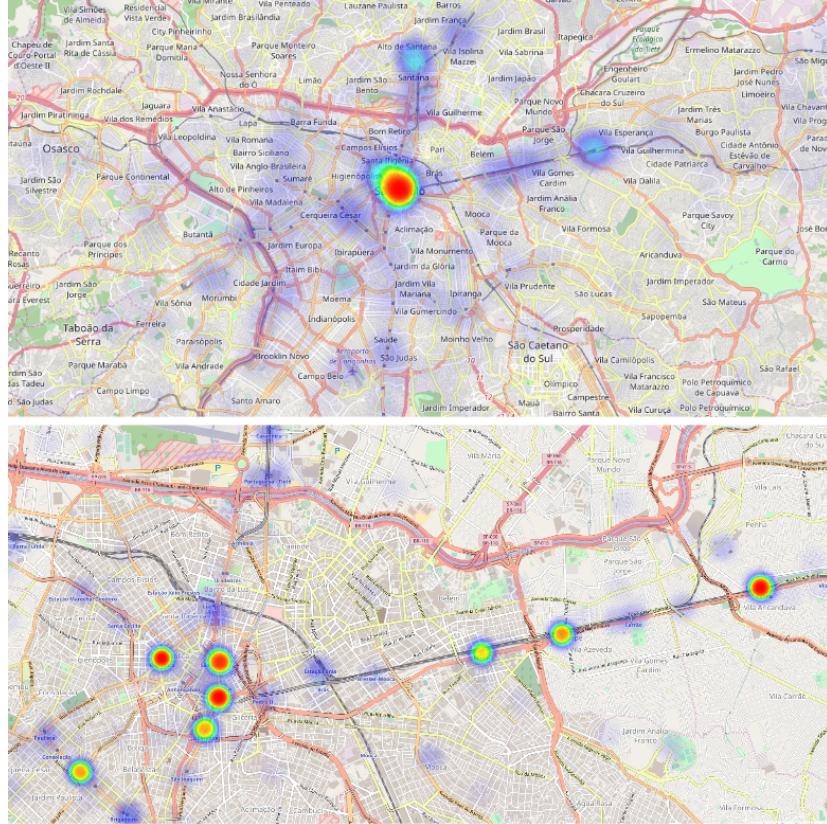


Figure 4.8: Rio de Janeiro Heatmap to the positive tweets



Experiments

Figure 4.9: São Paulo Heatmap to the positive tweets



In order to understand the spatial distribution of travel-related tweets we generated a heatmap
 1042 for both cities. From the heatmap of RJ, illustrated in Fig. ??, it is possible to identify that some
 agglomerations of tweets are located at Central do Brasil, Cidade Nova and Triagem train stations,
 1044 as well as at Uruguiana, Maracanã and Carioca metro stations. The Rio-Niterói bridge, connecting
 Rio de Janeiro to Niterói, as well as the piers on both sides also presented considerable clouds
 1046 of tweets classified as travel-related.

The heatmap for the city of SP, illustrated in Fig. ??, was also an interesting case to observe.
 1048 Almost every agglomeration matched some metro or train station. Estação Brás, Tatuapé, Belém,
 Estação Paulista, Sé, Liberdade were some of the stations highlighted in the heatmap. We could
 1050 also identify a little agglomeration of travel-related tweets at Congonhas airport, even though no
 tweets seemed to mention the word *plane* explicitly in the training of our classification model.

1052 **4.3.6 Final Remarks**

4.4 Travel-mode Extraction

Experiments

¹⁰⁵⁴ **Chapter 5**

Conclusions and Future Work

¹⁰⁵⁶ This planning report had two distinct objectives. The first one is the search of related works in
¹⁰⁵⁸ order to see what is already developed to the problem context of this dissertation. The second
¹⁰⁶⁰ objective is the initial planning of the dissertation work, as well as the approach and methodology
chosen to tackle the problem in hands. From the all work made so far, it is possible to make some
conclusions.

¹⁰⁶² This dissertation proposes to tackle the problem of extraction of aspect-based sentiment from
the citizens opinions about the services of a city, in social media streams, through a framework
that may be capable of processing the messages and build some appealing visual indicators.

¹⁰⁶⁴ Hence, the problem was decomposed in some sub-problems. The literature review served to
find interesting solutions for each sub-problem. There, a great diversity of approaches was found,
¹⁰⁶⁶ not only about sentiment analysis that is the most important task in this dissertation but also for
another problems like the content filtering and disambiguation.

¹⁰⁶⁸ The proposed framework can be seen as a potential tool to the users of the city's services and
for the responsible entities, allowing that only good decisions are made to improve the quality of
¹⁰⁷⁰ the cities and, in this particular case, the urban transportation systems.

¹⁰⁷² To summarize the conclusions of all the work made so far, a SWOT analysis was conceived
and the points that composed it are present below.

1. Strengths

- ¹⁰⁷⁴
 - Added value proposal by combining multiple State-of-the-Art approaches to tackle
chained sub-tasks;
 - Well defined sub-tasks/modules will make it easier to track errors.

2. Weaknesses

- ¹⁰⁷⁸
 - It might be difficult to collect enough relevant data for specific scenarios (e.g. the
quality of the urban transportation in Porto);
 - Twitter data might not be so reliable if there are few relevant messages.

Opportunities

- 1082 • New scenario application for aspect-based sentiment analysis: transportation systems
and Smart Cities;
- 1084 • Extending State-of-the-Art approaches in each sub-task/module if the target scenario
presents specific constraints.

Threats

- 1086 • Absence of ground-truths for the target scenarios may lead to underperformed mod-
ules;
- 1088 • Limited time for implementation is a risk of some unforeseen difficulties arise.

5.1 Expected Contributions

The work to be developed in this dissertation should present contributions both at the technological
1092 and scientific level. Some of the most important contributions are listed below:

- 1094 • A brief review of related literature to help contextualize readers in the subject of information
extraction, in particular the sentiment analysis, from social media streams and how difficult
is this task;
- 1096 • Development of a tool that could bring a potential value to the cities in order to improve the
quality of its services;
- 1098 • The studies of use cases about Smart Cities and Transportation Systems using aspect-based
sentiment analysis may be considered something innovative since there are very few works
1100 related with both scenarios.

5.2 Task Planning and Scheduling

1102 The tasks to be undertaken are mostly based in the modules described in Section ?? for the pro-
posed framework architecture. The first task is to choose what are the specific scenarios that will
1104 serve to test the developed framework. A priori, two different scenarios will be enough to prove
the good functionality and usability of the tool. Hence, the crawler module will be used to col-
1106 lect social media streams from the middle of February until, approximately, the ending of May.
Meanwhile, the setup of the framework environment needs to be done. After this first step, the
1108 development of the modules will occur. The first module to tackle is the aspect-based sentiment
analysis and the sub-module of preprocessing. With the estimation of a possible margin of error,
1110 these tasks are ready to employment in the begin or middle of April. The target filtering module
will be developed, if everything is going as planned, between the beginning/middle of April until
1112 the middle of May. The remaining month of May will serve to work on the aggregation module

Conclusions and Future Work

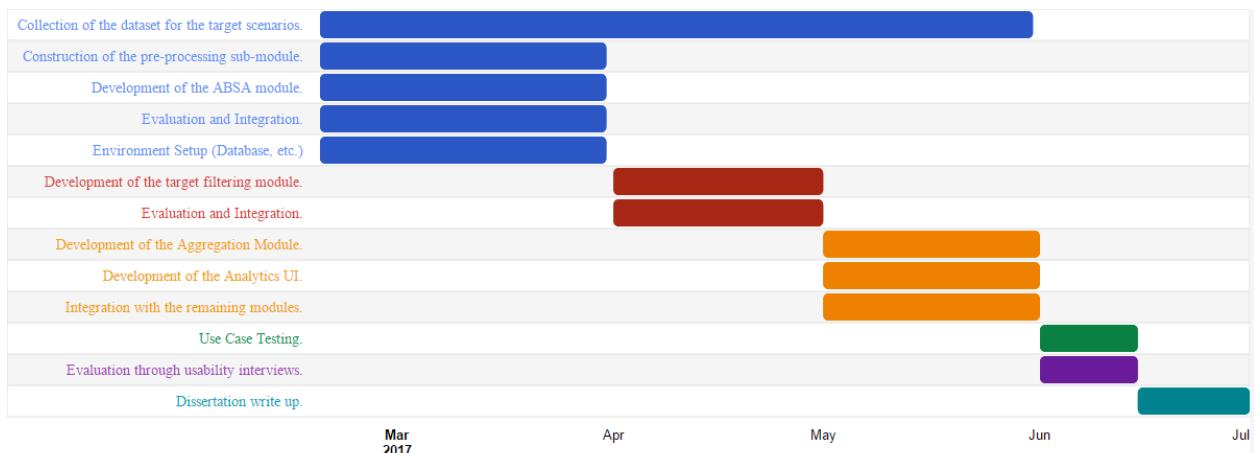


Figure 5.1: Dissertation working plan.

and the analytics UI. The month of June will be to test the final framework into the collected dataset about the two different scenarios. In order to evaluate the usability of the framework, it's planned the existence of a bunch of interviews to see if it's really possible that users of this tool are capable of immediately identify some conclusion from the analysis presented. This evaluation step will occur in the first two weeks of June, being the remaining two to the final dissertation report write up.

In the Figure ?? it's possible to visualize a Gantt chart scheduling according the mentioned tasks and the ideal scenario in case there are no delays.

Conclusions and Future Work

msjmnsm