

INSTITUTO TECNOLÓGICO DE AERONÁUTICA



Heládio Sampaio Lopes

**NATURAL LANGUAGE PROCESSING FOR TREND
FORECASTING**

Final Paper
2020

Course of Computer Engineering

Heládio Sampaio Lopes

**NATURAL LANGUAGE PROCESSING FOR TREND
FORECASTING**

Advisor

Prof. Dr. Filipe Alves Neto Verri (ITA)

COMPUTER ENGINEERING

SÃO JOSÉ DOS CAMPOS
INSTITUTO TECNOLÓGICO DE AERONÁUTICA

Cataloging-in Publication Data
Documentation and Information Division

Lopes, Heládio Sampaio
Natural Language Processing for Trend Forecasting / Heládio Sampaio Lopes.
São José dos Campos, 2020.
31f.

Final paper (Undergraduation study) – Course of Computer Engineering– Instituto Tecnológico de Aeronáutica, 2020. Advisor: Prof. Dr. Filipe Alves Neto Verri.

1. Natural Language Processing. 2. Deep Learning. 3. Machine Learning. I. Instituto Tecnológico de Aeronáutica. II. Title.

BIBLIOGRAPHIC REFERENCE

LOPES, Heládio Sampaio. **Natural Language Processing for Trend Forecasting**. 2020. 31f. Final paper (Undergraduation study) – Instituto Tecnológico de Aeronáutica, São José dos Campos.

CESSION OF RIGHTS

AUTHOR'S NAME: Heládio Sampaio Lopes

PUBLICATION TITLE: Natural Language Processing for Trend Forecasting.

PUBLICATION KIND/YEAR: Final paper (Undergraduation study) / 2020

It is granted to Instituto Tecnológico de Aeronáutica permission to reproduce copies of this final paper and to only loan or to sell copies for academic and scientific purposes. The author reserves other publication rights and no part of this final paper can be reproduced without the authorization of the author.

Heládio Sampaio Lopes
H8A St., 113
12228-460 – São José dos Campos–SP

NATURAL LANGUAGE PROCESSING FOR TREND FORECASTING

This publication was accepted like Final Work of Undergraduation Study

Heládio Sampaio Lopes

Author

Filipe Alves Neto Verri (ITA)

Advisor

Inaldo Capistrano Costa
Course Coordinator of Computer Engineering

São José dos Campos: JUNE 19, 2020.

Acknowledgments

Thank you

“That’s All Folks”
— Looney Tunes,

Resumo

Resumo

Abstract

Abstract

List of Figures

FIGURE 2.1 – Stemming process for “connect” variations, Figure from (VIJAYARANI <i>et al.</i> , 2015).	17
FIGURE 2.2 – Bag of Words example.	18
FIGURE 2.3 – Word2Vec architectures, Figure from (MIKOLOV <i>et al.</i> , 2013).	19
FIGURE 2.4 – Tri-gram representation for “apple” word.	20
FIGURE 3.1 – Flowchart of the proposed framework, Figure from (HURTADO <i>et al.</i> , 2016).	23
FIGURE 3.2 – Ensemble forecaster framework, Figure from (HURTADO <i>et al.</i> , 2016).	24
FIGURE 4.1 – Database time split representation.	27
FIGURE 4.2 – Topic identification process.	28
FIGURE 4.3 – Multi-class classifier from <i>Labeled</i> set.	28
FIGURE 4.4 – Flowchart to evaluate the time series model.	29

List of Tables

TABLE 5.1 – Tasks schedule over the months.	30
---	----

List of Abbreviations and Acronyms

AI	Artificial Intelligence
BoW	Bag of Words
CBOW	Continuous Bag of Words
NLP	Natural Language Processing
ML	Machine Learning
TF-IDF	Term Frequency Inverse Document Frequency
IT	Information Technology

Contents

1	INTRODUCTION	14
1.1	Motivation	14
1.2	Objective	14
1.3	Organization of this work	15
2	LITERATURE TO REVIEW	16
2.1	Natural Language Processing	16
2.1.1	Text Processing Techniques	16
2.1.2	Word Embedding	18
2.1.3	Topic Clustering	21
2.2	Machine Learning	21
3	RELATED WORKS	22
3.1	Topics Discovery	22
3.2	Trend Forecast	23
3.3	Final Remarks	25
4	MATERIALS AND METHODS	26
4.1	Objective	26
4.2	Research method	26
4.2.1	Database	27
4.2.2	Pre-processing the data	27
4.2.3	Topic identification	27
4.2.4	Document classification	28

4.2.5	Forecast evaluation	29
5	ROADMAP	30
	BIBLIOGRAPHY	31

1 Introduction

Recently, the growth of artificial intelligence has been helping us to solve problems in the most various areas, including in linguistics. With natural language processing techniques, computers are able to process text in order to extract information faster than humans.

1.1 Motivation

Every kind of expression, verbal or in writing, brings us a lot of information to be interpreted. Whether the topic is chosen, the tone used or the choice of words, everything can be interpreted, and then generate some useful information. Over the years, more and more knowledge is generated and we humans are not able to process such an amount of information. Natural language processing emerges as a technology capable of assisting us in this hard task.

Khurana *et al.* (2017) defines natural language processing, abbreviated by NLP, as a branch of artificial intelligence capable of making computers understand and extract information from human language. NLP can perform a lot of tasks, such as identifying different topics for a set of documents, classifying texts on predefined subjects, and beyond that extract the sentiment to know what people are saying about something.

1.2 Objective

Curious about the fast world's evolution, this work aims to explore and compare several Natural Language Processing techniques to model the topic's evolution over time. With this in mind, evaluate the ability of those models to make predictions about future trends.

1.3 Organization of this work

The remaining of this work is organized as follows: Chapter 2 will cover the theory behind Natural Language Processing. Chapter 3 will describe some previous works which use topic discovery and trend forecast. Chapter 4 will better explain the problem and the methodology that will be used to treat it. Finally, Chapter 5 will present the chronological roadmap until the end of the work.

2 Literature to Review

In this chapter, we will introduce the general concepts and techniques behind Natural Language Processing. We will cover all the necessary steps for extracting meaningful topics from texts.

2.1 Natural Language Processing

2.1.1 Text Processing Techniques

The key task to several machine learning problems consist in make a good data processing before applying any model. A clean data set can allow a model to increase its performance in the learning process, making a better identification in the patterns present in the variables. Hence, in the next sections, it will be discussed a few techniques to clear the text and prepare it for ML algorithms.

2.1.1.1 Normalization

There is no right way to normalize text, this process has it is really important to put all text in the same level. A normalization process has a series of steps to be followed sequentially, all of them can be seen as 4 big tasks: stemming, lemmatization, stop words removal and everything else.

1. Stemming: Is the process of reduce inflected words to a primitive form, the stem. This method is able to remove the word's affixes to capture its base meaning, and still reducing the number of variations to save memory space. Figure 2.1 shows how some inflections for “connect” can be converted to its root form.
2. Lemmatization: similar to stemming, this step also reduce words to some primitive form, but with a little improvement. Lemmatization can returns the words to his dictionary form, based on its part of speech context. So it is possible to discriminate words with the same spelling but different meanings depending on the context.

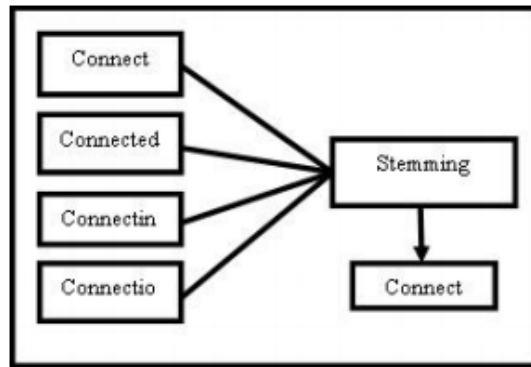


FIGURE 2.1 – Stemming process for “connect” variations, Figure from (VIJAYARANI *et al.*, 2015).

3. Remove stop words: Many words can occurs a several times in a document without add any meaningful information, such as *the*, *is*, *at*, *which*, and *on*. Their high frequency can be seen as an obstacle to perform good results on NLP models, (KANNAN; GURUSAMY, 2014).

There are some types to remove stop words, most of then based on evaluating the frequency of words in text, for more information see (VIJAYARANI *et al.*, 2015). But the classic and easier method is based on using a pre-compiled list of know words and removing then from text.

4. Everything else: Differently from the previous steps, the last one doesn’t need any grammar rules or even a frequency analysis, it’s purely text manipulation. It involves set all character to lowercase; remove numbers or convert then to word form; remove punctuation; expand contractions; convert special characters to ASCII form; and any other conversion needed.

2.1.1.2 Tokenization

Once the data is normalized, we need to know how to represent it. The tokenization process consists in splitting longer strings into meaningful small pieces called tokens. The most common way to tokenize a text is chunking it the into words, ie, given a piece of text the tokenize process will return a list of words.

2.1.1.3 Bag of Words

The machine learning algorithms take numerical features as input, hence, it will bee necessary to represent the text in numerical form. With the Bag of Words model we can represent in matrix form a set of documents.

With the tokenization output we will have the lists representations for all documents

in the data set. Those lists can be interpreted as vectors over the vector space of all unique tokens, also called by vocabulary. So, for a given sentence, we mark how many times its words appears in the list indexes where each entry corresponds to a word in the vocabulary. The Figure 2.2 show a simple example of how three sentences can be represented with BoW model.

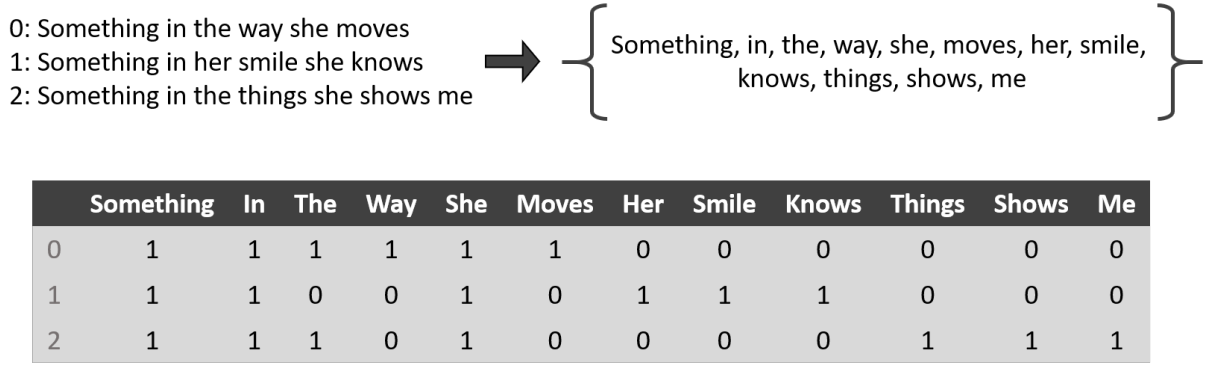


FIGURE 2.2 – Bag of Words example.

2.1.1.4 TF-IDF

Term Frequency Inverse Document Frequency, TF-IDF for short, it is applied to a BoW to determine the relative frequency for words in a specific document when compared to the inverse proportion of that word over all documents in the collection. So, it can be determined how important are the words in a specific document.

From BoW, for the i^{th} vocabulary's word in the j^{th} document, its TF-IDF weight is:

$$w_{i,j} = \text{tf}_{i,j} \times \log \left(\frac{N}{\text{df}_i} \right). \quad (2.1)$$

Where, the term frequency, $\text{tf}_{i,j}$, is how many time i^{th} word appears in the j^{th} document. The document frequency, df_i , is the number of documents in which th i^{th} vocabulary words is present. And, finally, N is the size of the document collection, with a large number of documents this term can explodes, so the logarithmic function is applied to dampen this effect.

2.1.2 Word Embedding

The vectorization methods such as BoW and TF-IDF can be very useful, but they can not represent the words context. This means that the same words used in different contexts have the same representation, just as different words used with the same meaning

are represented differently. Besides that, an one-hot encoding method, like BoW, presents a very sparse representation with high dimensionality.

The Word Embedding is a technique to represent words in vectors capable of capture the words context in a document. It is also able to smooth the high dimensionality effect by using much more compact vector to represent the words.

There are three most know way to perform a good word embedding. We will describe briefly each one of them below.

2.1.2.1 Word Representations in Vector Space

The first great word embedding technique emerged when Google researchers proposed two architectures to build continuous vector representations of words. Word's context can be observed as the words that surround it in a sentence. Then using shallow neural networks, it is possible to calculate the word vector space based on word's context, (MIKOLOV *et al.*, 2013).

The first suggested approach is the continuous bag of words or CBOW, the left side of Figure 2.3 show its architecture. Here the neural networks is designed to predict, given the context, which word is most likely to appear. So, words with the same probability to appear can have a shared dimension in the words vector space.

The second approach is known by Skip-Gram, architecture at right in Figure 2.3. Very similar to CBOW, but instead of predicting the current word the Skip-Gram uses the current word as an input to a neural network to predict its context.

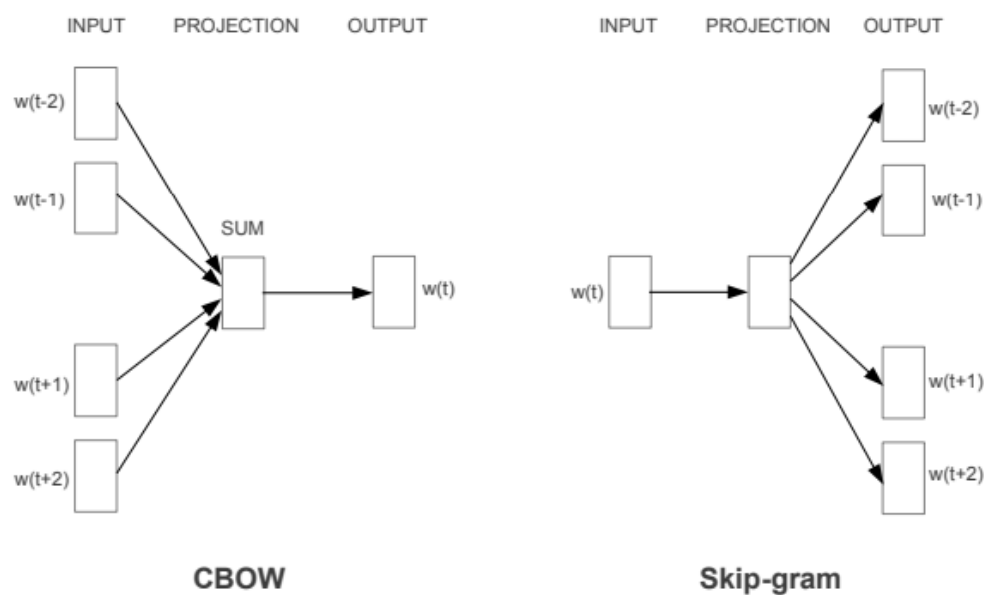


FIGURE 2.3 – Word2Vec architectures, Figure from (MIKOLOV *et al.*, 2013).

After the network training process we can use the hidden layer weight matrix as an lookup table to build the word embedding representation. The dimension for the vector space is managed by the number of neurons in the hidden layer.

2.1.2.2 Global Vectors for Word Representation

Just a year later Pennington *et al.* (2014) arrives with a new approach to represent words in a vector space. The Global Vectors for Word Representation, or GloVe, method emerged by the need to consider some factors ignored by Skip-Gram.

Methods such as Skip-Gram learn their embedding by targeting words to their respective context, ignoring the fact that some words appear more in a context than others. Thus, this co-occurrence of words only adds more useless training examples, increasing the training complexity without adding relevant information.

GloVe, however, proposes to use the corpus statistics in a more efficient way. Using a weighted least squares model trained on a global word-word co-occurrence counts matrix. Thereby, it is possible to build a lookup table for the words in vocabulary and use it to represent them in a vector space.

2.1.2.3 Word Vectors with Subword Information

Both Skip-Gram and GloVe provide a good vector representation for words, but there still are an unsolved problem, What to do with unknown words? To solve this question was proposed a new embedding technique which uses subword units to build a vector space, (BOJANOWSKI *et al.*, 2017).

Similar to Skip-Gram, this new method, the FastText, train its embedding by using a target to predict the context. However, instead of using the full words FastText goes a level deeper, breaking the words in n -grams, ie, the word becomes its own context. The Figure 2.4 shows how the word “apple” can be broken into n -grams.

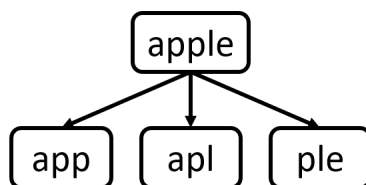


FIGURE 2.4 – Tri-gram representation for “apple” word.

There are a couple great advantages by using this method. It is now possible to generalize new words, or unseen in training data, since they have the same characters as known ones. Although it is possible to use available pre-trained models, the FastText

requires less text to be trained, it can extract much more information from small pieces of text.

2.1.3 Topic Clustering

2.2 Machine Learning

3 Related Works

In the last chapter, we saw the theoretical foundation on NLP techniques. In this chapter, we will review in the literature some works that use the NLP techniques described to discover topics in a data set. In addition, we will show some applications for this type of task. And, finally, some final remarks to continue this work.

3.1 Topics Discovery

Finding meaningful topics in a document collection has been used for a lot of authors for the most various applications. For example, Hurtado *et al.* (2016) use topic modeling to inspect research publications, patents, and technical reports aiming to model the evolution of the direction of research and forecast the near future trends in IT industry.

Using the titles and abstracts of a data set with more than six thousand academic papers between 2002 and 2010, mostly collected by Tang *et al.* (2008), they proposed a sentence-level association rule to discover the meaningful topics. After categorizing the documents in topics, they were capable of building time series for each found topic, marking how many times that topic was cited in a given year. So, they were able to build an ensemble of forecasters to study the patterns and relationships among topics over the years.

For a better understanding, the Figure 3.1 has a flowchart with their proposed framework for the topic discovery and forecasting.

This framework involves some well-known major steps of NLP processing. First, they convert the documents into a transactional form, i.e., the phrases in each document will be considered individually during the process. Next, they perform the basic normalization steps which includes case conversion, tokenization, removing stop words, part of speech tagging, stemming and lemmatization. It is also performed an additional step, specific to their application, removing verbs such as “exploiting”, “adapting” and “propose”, because they are very common in scientific publications and do not add much meaningful information.

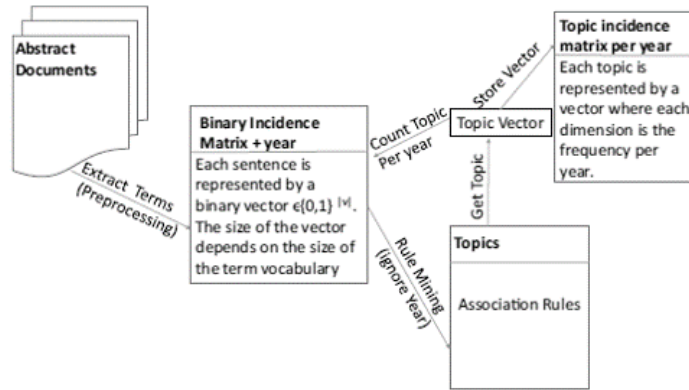


FIGURE 3.1 – Flowchart of the proposed framework, Figure from (HURTADO *et al.*, 2016).

To vectorize the transactions, it is used a slight variation of BoW. Instead of word counting, it is only checked whether a word belongs to a transaction, this is called the binary incidence matrix. The topic discovery step comes afterwards, applying an association rule mining to the transactions and discovery their patterns. In order to avoid different topics with redundant words, is applied a rule refinement process that allows similar topics to be combined.

Online forums and social media are excellent platforms for people discuss and share information about the most kind of subjects. Recently, the topic discovery technique was used to summarize different topics related to COVID-19 disease and perform a sentiment analysis on them (JELODAR *et al.*, 2020).

Reddit is a discussion website in which its users can submit posts and start discussions with other community members. The posts are organized in the called “sub reddits”, boards created by users to discuss a specific subject. Using over half million comments from 10 health related sub reddits with information about COVID-19, Jelodar *et al.* (2020) performed a topic discovery to group similar comments. So, applying a sentiment analysis on each comment it was possible to summarize the average opinion about the discovered topics.

3.2 Trend Forecast

Predicting future trends can be very helpful in various applications, like to model the evolution of research. Following the topic discovery process from Hurtado *et al.* (2016), a forecast trend was used to predict the near future related to each discovered topic. With all documents belonging to at least one identified topic in the set, was created a topic incidence matrix that contains the count of times a topic is mentioned over the years. Finally, they make a ensemble forecasting to predict the future topic counting using the

framework shown in Figure 3.2.

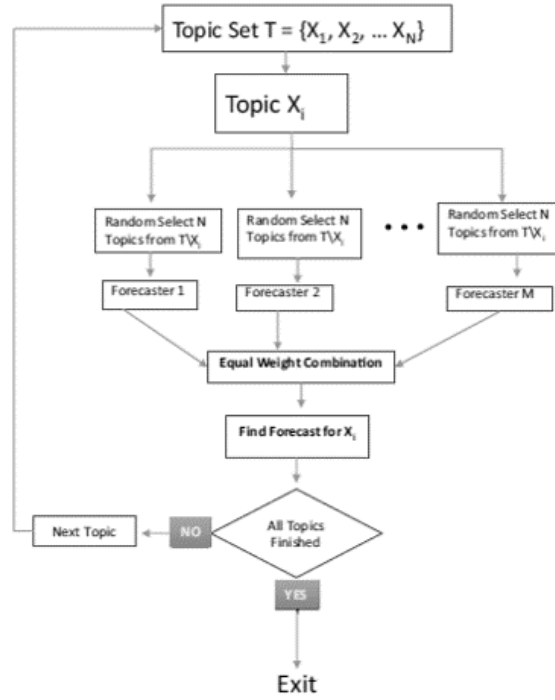


FIGURE 3.2 – Ensemble forecaster framework, Figure from (HURTADO *et al.*, 2016).

Given a specific topic, they generate M forecasters which target X_i along with N randomly chosen fields, excluding X_i . Then, the predicted value, $\hat{X}_i(t+1)$, is an average of each individual forecast, $\hat{X}_{i,F_k}(t+1)$, calculated by

$$\hat{X}_i(t+1) = \frac{1}{M} \sum_{k=1}^M \hat{X}_{i,F_k}(t+1). \quad (3.1)$$

By evaluating metrics like the coefficient of determination (R-squared) and mean squared error (MSE) they were able to conclude the ideal number N of variables to predict more accurately the future for all topics in the set.

Shen (2018) also predicted trends by analyzing the exponential growth in the volume of scholarly articles published over the years. However, he skips the NLP process step by using a pre-labeled data set from Springer containing the number of works above 14 subjects in 25 years. Using ensemble forecast based on neural networks and support vector regression he was able to study the topics growth and codependency between them.

3.3 Final Remarks

As we saw earlier topic discovery and trend forecast were subjects widely explored in the literature. With that in mind, we wish, in this work, to reproduce these techniques. However, in addition to what has been presented we want to be able to explore some modifications.

Assuming a system that receives documents in real-time, let's assume they are news, it would be very computationally expensive to redo all the topic discovery process with each new news. It would be much simpler if there was a process that, given an input document, would be able to identify which topics were covered by it.

So, aiming at this type of application, we will propose a system capable of such activity to model the topics evolution over the time.

4 Materials and Methods

In this chapter, we will discuss the roadmap steps to carry out this work. First, we will reintroduce the work objectives. Finally, the work plan with the necessary steps to accomplish the proposed objectives.

4.1 Objective

As already discussed, we want to build models capable of making predictions regarding the evolution of discovered topics in a set of documents. We also want to find topics in real time at each received document without having to redo the topic discovery process. Then, by the end of this work we must have performed the tasks listed below.

- Find a database long enough, over several years;
- Perform all necessary treatment steps to normalize the documents' texts;
- Find meaningful topics in a subset from the original database;
- Create a topic classifier to find out if a document addresses any topic of interest;
- Model topic evolution to evaluate the forecast accuracy.

4.2 Research method

With the objectives defined, we must have action plans to achieve them. Next, we will discuss in more detail the plan action for each of these previous punctuated steps.

We need to mention the figure notation used in next sections. Whenever a box with the bottom right tip folded appears, it represents a machine learning process with all the applicable steps, such as cross-validation and hyperparameter tuning.

with all cross-validation steps and tuning of applicable parameters

We need to mention the math notation used in this work. The neuron a_{ji} represents the i th neuron in the j th layer of network.

4.2.1 Database

The first necessary task is to find a database over several years. Some options are available, such as daily news from wikipedia, newspapers articles, research and academic papers, patents, and even social media data like twitter or reddit. We must be careful to choose a database which covers several years in order to get a good temporal representation when the modeling their topics.

4.2.2 Pre-processing the data

With the chosen database, we must define a data pre process pipeline to normalize the documents, putting them all at the same pattern. For this, we can use the normalization techniques shown in Chapter 2 like stemming, lemmatization, stop words removal and all necessary textual manipulations to obtain the best text representation for the documents.

After processing the data we need to index them over the time and, then, split the full treated data set in three subsets. They must be time ordered, as shown in Figure 4.1, this means that for each document in T_i its time index $t_x^{(i)}$ must be lower than any index $t_x^{(m)}$ in a document contained in T_m , and so on.

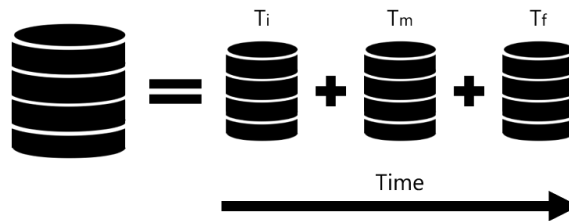


FIGURE 4.1 – Database time split representation.

Before proceeding, let's name the subsets. The initial subset will be called *Identifier* from now on, the middle and final ones will be respectively designated by *Modeler* and *Validation*.

4.2.3 Topic identification

With the *Identifier* set, techniques of document clustering will be applied to identify the discussed subjects in the documents, the Figure 4.2 illustrates the steps sequence

of this task. Similar to Hurtado *et al.* (2016), a refinement will be made so that only significant topics remain.

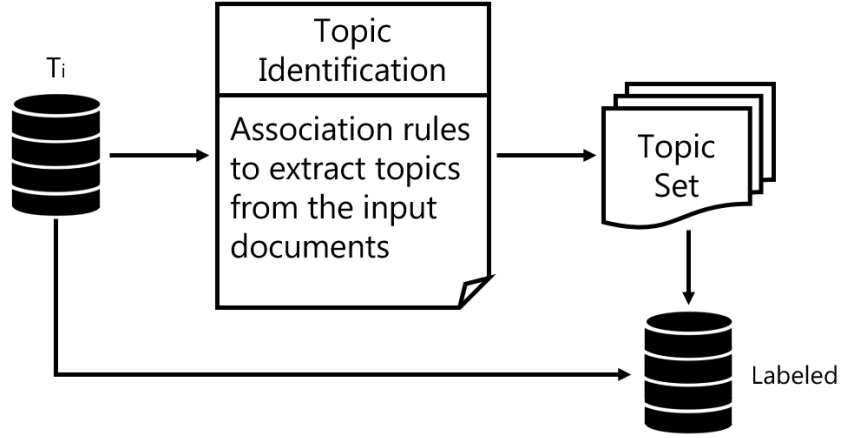


FIGURE 4.2 – Topic identification process.

Next to the topics identification we will have a topic set, then each document contained in *Identifier* can be labeled with at least one topic from the set, this new database will be called by labeled *Identifier*.

4.2.4 Document classification

For each topic in our set of discovered topics, we must be able to identify which topics are covered by a new document. Thus, we will build a classifier to perform this verification.

Knowing that a document can talk about several topics, so we must have a multi-class classifier. We can see this as an individual binary classifier for each topic that tells us whether the document has it. Using the labeled *Identifier* set it is possible to build this classifier. Figure 4.3 illustrates it in detail.

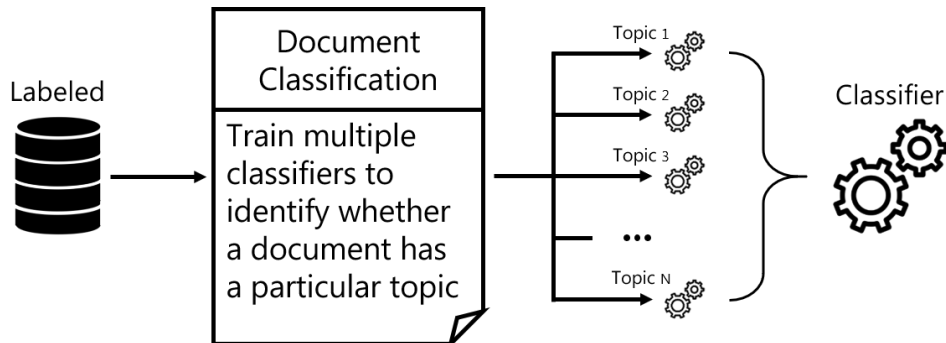


FIGURE 4.3 – Multi-class classifier from *Labeled* set.

4.2.5 Forecast evaluation

Finally, evaluate the time series model is the last task to be accomplished. Following the flowchart shown in Figure 4.4, first, we have to label the *Modeler* and *Validation* sets. Then, using labeled *Identifier* and *Modeler* we will build a topic incidence matrix over the time, to apply a forecaster process for those time series. With the labeled *Validation* set we will perform an evaluation for our model and then make conclusions about it.

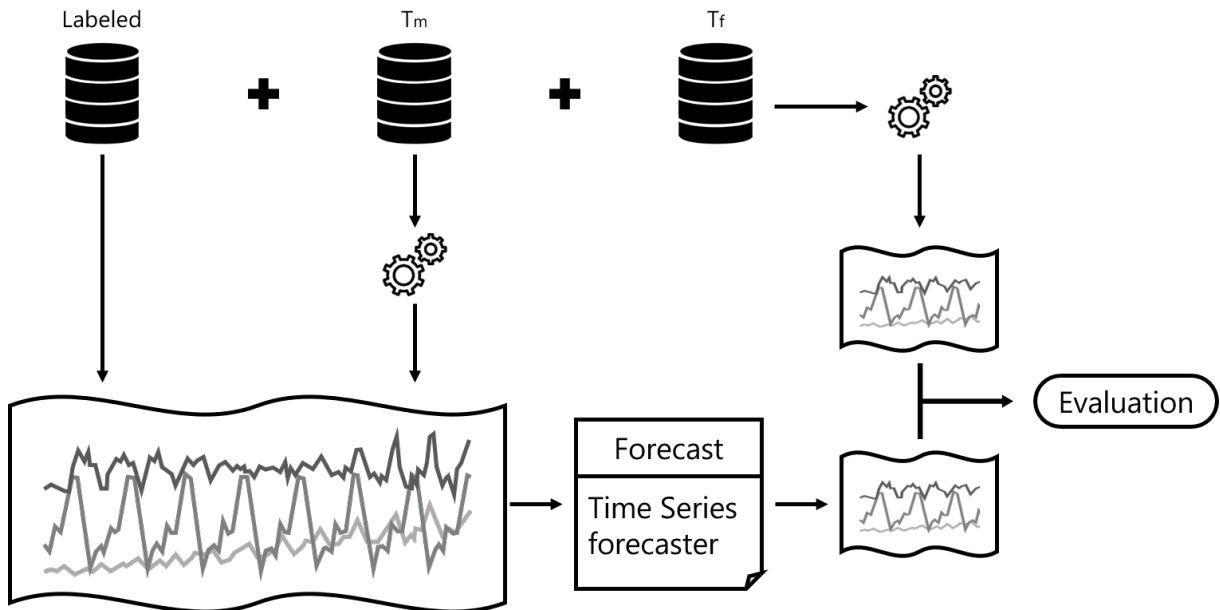


FIGURE 4.4 – Flowchart to evaluate the time series model.

5 Roadmap

In view of the problem’s complexity, we can elaborate a schedule with the proposed tasks in the previous chapter. The Table 5.1 show the tasks over the remains months until the end of this work.

TABLE 5.1 – Tasks schedule over the months.

Month	Task	Description
July		
August		
September		
October		
November		

Bibliography

BOJANOWSKI, P.; GRAVE, E.; JOULIN, A.; MIKOLOV, T. Enriching word vectors with subword information. **Transactions of the Association for Computational Linguistics**, v. 5, p. 135–146, 2017. ISSN 2307-387X.

HURTADO, J. L.; AGARWAL, A.; ZHU, X. Topic discovery and future trend forecasting for texts. **Journal of Big Data**, Springer, v. 3, n. 1, p. 7, 2016.

JELODAR, H.; WANG, Y.; ORJI, R.; HUANG, H. Deep sentiment classification and topic discovery on novel coronavirus or covid-19 online discussions: Nlp using lstm recurrent neural network approach. **arXiv preprint arXiv:2004.11695**, 2020.

KANNAN, S.; GURUSAMY, V. Preprocessing techniques for text mining. **International Journal of Computer Science & Communication Networks**, v. 5, n. 1, p. 7–16, 2014.

KHURANA, D.; KOLI, A.; KHATTER, K.; SINGH, S. Natural language processing: State of the art, current trends and challenges. 08 2017.

MIKOLOV, T.; CHEN, K.; CORRADO, G.; DEAN, J. Efficient estimation of word representations in vector space. **arXiv preprint arXiv:1301.3781**, 2013.

PENNINGTON, J.; SOCHER, R.; MANNING, C. Glove: Global vectors for word representation. In: **Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)**. Doha, Qatar: Association for Computational Linguistics, 2014. p. 1532–1543. Disponível em: <<https://www.aclweb.org/anthology/D14-1162>>.

SHEN, Q. Topic discovery and future trend prediction in scholarly networks. In: . Shanghai, China: [s.n.], 2018.

TANG, J.; ZHANG, J.; YAO, L.; LI, J.; ZHANG, L.; SU, Z. Arnetminer: extraction and mining of academic social networks. In: **Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining**. [S.l.: s.n.], 2008. p. 990–998.

VIJAYARANI, S.; ILAMATHI, M. J.; NITHYA, M. Preprocessing techniques for text mining-an overview. **International Journal of Computer Science & Communication Networks**, v. 5, n. 1, p. 7–16, 2015.

FOLHA DE REGISTRO DO DOCUMENTO

1. CLASSIFICAÇÃO/TIPO TC	2. DATA June 19th, 2020	3. DOCUMENTO N° DCTA/ITA/DM-018/2015	4. N° DE PÁGINAS 31
5. TÍTULO E SUBTÍTULO: Natural Language Processing for Trend Forecasting			
6. AUTOR(ES): Heládio Sampaio Lopes			
7. INSTITUIÇÃO(ÕES)/ÓRGÃO(S) INTERNO(S)/DIVISÃO(ÕES): Aeronautics Institute of Technology – ITA			
8. PALAVRAS-CHAVE SUGERIDAS PELO AUTOR: Natural Language Processing; Deep Learning; Machine Learning			
9. PALAVRAS-CHAVE RESULTANTES DE INDEXAÇÃO: Natural Language Processing; Deep Learning; Machine Learning			
10. APRESENTAÇÃO:		(X) Nacional () Internacional	
ITA, São José dos Campos, 2020. Trabalho de Graduação. 31 páginas.			
11. RESUMO: Resumo			
12. GRAU DE SIGILO: (X) OSTENSIVO () RESERVADO () SECRETO			