

# Configuration Manual

MSc Research Project Data Analytics

# Priscila Cristina da Silva de Oliveira x23157003

School of Computing National College of Ireland

Supervisor: Barry Haycock

#### **National College of Ireland**



#### **MSc Project Submission Sheet**

#### **School of Computing**

Student Name:	Priscila Cristina da Silva de Olive	ira

**Student ID:** x23157003

**Programme:** MSc in Data Analytics

**Year:** 2024

**Module:** Research Project

Lecturer: Bubmission Due

Barry Haycock

**Date:** 16.09.2024

**Project Title:** From LDA to BERTopic: Evaluating Topic Modelling Methods for

Aviation Safety Reports in Brazilian Portuguese

Word Count: 3149 Page Count: 17

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Privala C. D. alnera
Date:	16.09.2024

#### PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple	
copies)	
Attach a Moodle submission receipt of the online project	
<b>submission,</b> to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project, both	
for your own reference and in case a project is lost or mislaid. It is not	
sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only		
	Office Use Only	

Signature:	
Date:	
Penalty Applied (if applicable):	

# Configuration Manual

# Priscila Cristina da Silva de Oliveira x23157003

# 1 System Specifications

This section details the hardware and operating system used in this study. The hardware specifications are not minimum requirements. However, one should be aware that different specifications could result on different performance.

#### 1.1 Hardware Specifications

Processor	13th Gen Intel(R) Core(TM) i9-13980HX 2.20 GHz
RAM	64.0 GB
GPU	NVIDIA GeForce RTX 4070 Laptop GPU

#### 1.2 Software Environment

Operating System	Windows 11 Pro, 64-bit
Integrated Development Environment	PyCharm 2024.1.1
Jupyter Notebook	6.5.4
Python Version	3.11.7

#### 1.2.1 Key libraries

Library / Package	Version
adjustText	1.2.0
beautifulsoup4	4.12.2
bertopic	0.16.3
gensim	4.3.0
matplotlib	3.8.0
nltk	3.8.1
numpy	1.26.4
pandas	2.1.4
pyLDAvis	3.4.1
PyMuPDF	1.24.7
seaborn	0.12.2
scikit-learn	1.2.2
scipy	1.12.0
resquests	2.31.0

tqdm	4.65.0
umap-learn	0.5.6
wordcloud	1.9.3

#### 2 Dataset Characteristics

The aviation safety reports used for the research project are publicly available on the website <a href="https://sistema.cenipa.fab.mil.br/cenipa/paginas/relatorios/relatorios.php">https://sistema.cenipa.fab.mil.br/cenipa/paginas/relatorios/relatorios.php</a> in the PDF format. The documents are distributed over 31 tabs and most of the files are in Portuguese. Figure 1 shows a snippet of the webpage. The flags represent the language in which the report is available.

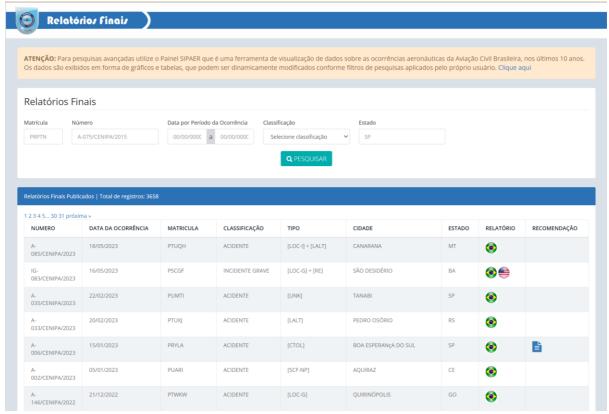


Figure 1 Reports in CENIPA website

Some documents are available in English and Spanish. It was also observed that some of the files, especially those reporting older events, are scanned documents which are not machine-readable. In addition, during the download phase, it was noted that a small number of documents cannot be accessed although the instance is shown in the reports table.

The download of the files was performed using web scraping leveraging BeautifulSoup4. The first version of the code used did not consider the different existing tabs and, for this reason, only the documents in the first tab were downloaded. The steps followed are described in Algorithm 1. The code described in the algorithm is available in the file named download\_pdf\_version\_1.ipynb.

#### **Algorithm 1** Download PDF files from CENIPA website – version 1

```
FUNCTION main():
```

```
SET URL to "https://sistema.cenipa.fab.mil.br/cenipa/paginas/relatorios/relatorios.php"
```

response = GET request to URL

IF response status is not successful THEN PRINT error message RETURN false

page = Parse HTML content of response

table = Find element in page with ID "lista"

links = Extract all <a> elements from table

FOR EACH link in links:

href = Get "href" attribute of link

IF "pdf" is not in href THEN
CONTINUE to next iteration

filename = Remove path from href, keeping only the file name

pdfUrl = Concatenate "https://sistema.cenipa.fab.mil.br/cenipa/paginas/relatorios/" and href

pdfContent = GET request to pdfUrl

Write pdfContent to file named filename

PRINT "Getting {filename}..."

RETURN true

CALL main()

Given that not all the documents available are machine-readable and that this project focuses on the reports written in the Portuguese language, the initial code was modified to:

- Classify the reports between "Portuguese" and "Other Languages"
- Subclassify the reports between "Scanned" and "Searchable"
- Access the 31 existing tabs to download all the files available.

The final code for the web scrapping process verifies whether the file exists before proceeding with the download. If so, it skips to the next file. It also informs in cases when it fails to download the document.

In order to verify whether a certain file is machine-readable, the function *is\_searchable()* was defined. The *main()* function defines the folder and subfolder names. It verifies if the main folders exist, before creating them.

Next, it iterates over the 31 existing tabs using *BeautifulSoup* to parse the HTML content. If the PDF link is valid, it categorises the file and save or skip the document. Algorithm 2 describes the steps used to download the files used in the study and the code can be found in the file download\_pdf\_final\_version.ipynb.

#### Algorithm 2 Download PDF files from CENIPA website – final version

FUNCTION is\_searchable(pdf\_path):

TRY:

Open PDF at pdf\_path
IF PDF has no pages THEN
RETURN false

Get text from first page RETURN true if text is not empty, false otherwise

CATCH IndexError:

PRINT error message RETURN false

FUNCTION main():

SET portuguese\_folder to 'PortugueseReports'
SET other\_languages\_folder to 'Other\_Languages'
SET subfolders to ['Scanned', 'Searchable']

CREATE folders and subfolders if they don't exist

FOR page number FROM 1 TO 31:

SET URL to

 $f"https://sistema.cenipa.fab.mil.br/cenipa/paginas/relatorios/relatorios?\&\&?\&pag=\{page number\}"$ 

Send GET request to URL IF request failed THEN PRINT error message RETURN false

Parse HTML content of response Find table with id "lista" Extract all links from table

FOR EACH link in links:

IF link doesn't contain 'pdf' THEN CONTINUE to next link

SET pdf\_url to full URL of PDF Download PDF file

IF download successful THEN
Save PDF to temporary file

IF link title is "Relatório Final em Português" THEN SET language\_folder to portuguese\_folder ELSE

SET language\_folder to other\_languages\_folder

IF is\_searchable(temporary\_file) THEN

SET save\_subfolder to 'Searchable'

**ELSE** 

SET save\_subfolder to 'Scanned'

SET save\_path to language\_folder/save\_subfolder/pdf\_filename

IF save\_path doesn't exist THEN

Move temporary file to save\_path

PRINT success message

**ELSE** 

Delete temporary file PRINT file already exists message

**ELSE** 

PRINT download failure message

RETURN true

CALL main()

# 3 Extracting Metadata – Types of Occurrences

The types of occurrences were manually extracted from CENIPA's website. The values were copied and pasted on Microsoft Excel. The unique values are shown in Table 1. The code ALTL was identifies among the types of occurrences. However, by analysing the report classified in the category, it became clear that it was a typing error. The occurrence should have been classified as LALT.

**Table 1: Types of Occurrences** 

Type	Description
ADRM	aerodrome
AMAN	abrupt maneuvre
ARC	abnormal runway contact
ATM/CNS	atm/cns
BIRD	bird
CFIT	controlled flight into/toward terrain
CTOL	collision with obstacle(s) during take-off and landing
EXTL	external load related occurrences
F-NI	fire/smoke (non-impact)
F-POST	fire/smoke (post-impact)
FUEL	fuel related
GCOL	ground collision
GTOW	glider towing related events
ICE	icing
LALT	low altitude operations
LOC-G	loss of control – ground

LOC-I	loss of control – inflight
LOLI	loss of lifting conditions en-route
MAC	airprox/tcas alert/loss of separation/near midair collisions/midair collisions
MED	medical
OTHR	other
RAMP	ground handling
RE	runway excurision
RI	runway incursion
SCF-NP	system/component failure or malfunction (non-powerplant)
SCF-PP	system/component failure or malfunction (powerplant)
TURB	turbulence encounter
UIMC	unintended flight in imc
UNK	unknown or undetermined
USOS	undershoot/overshoot
WILD	collision with animals
WSTRW	wind shear or thunderstorm

#### 4 LDA

Algorithm 3 describes the steps taken for the LDA model. The respective code is available in the file named LDA\_final\_version.ipynb.

#### Algorithm 3 Topic Modelling with LDA

SET seed to 23157003

INITIALIZE random number generators with seed

**DOWNLOAD NLTK resources** 

SET pdf\_folder\_path to location of PDF files

DEFINE extra\_stop\_words

COMBINE extra\_stop\_words with Portuguese stopwords from NLTK

FUNCTION extract\_text\_from\_pdf(pdf\_path):

INITIALIZE empty text string

TRY:

OPEN pdf\_file

FOR EACH page in pdf\_file:

ADD page text to text string

CATCH any exceptions:

PRINT error message

RETURN text string

FUNCTION preprocess(text):

CONVERT text to lowercase

REMOVE numbers from text

REMOVE punctuation from text

TOKENIZE text

RETURN list of tokens not in stop\_words

 $FUNCTION\ get\_top\_ngrams(processed\_texts,\ n,\ top\_k):$ 

CREATE empty list for n-grams

FOR EACH text in processed\_texts:

GENERATE n-grams from text

ADD n-grams to list

COUNT frequency of n-grams

RETURN top\_k most common n-grams

FUNCTION compute\_coherence\_and\_perplexity(dictionary, corpus, texts, start, limit, step):

INITIALIZE empty lists for coherence\_values, perplexity\_values, and model\_list

FOR num\_topics FROM start TO limit STEP step:

CREATE LdaModel with num\_topics

COMPUTE coherence score

**COMPUTE** perplexity

ADD scores and model to respective lists

RETURN model\_list, coherence\_values, perplexity\_values

FUNCTION evaluate\_lda\_model(corpus, dictionary, processed\_texts, start, limit, step):

CALL compute\_coherence\_and\_perplexity

NORMALIZE coherence and perplexity scores

COMPUTE composite scores

FIND best model based on highest composite score

PRINT best model details

PLOT evaluation metrics

RETURN best\_model, best\_num\_topics, best\_coherence, best\_perplexity

FUNCTION plot\_evaluation\_metrics(start, limit, step, coherence\_values, perplexity\_values, composite\_scores):

CREATE plot with three y-axes

PLOT coherence scores

PLOT perplexity scores

PLOT composite scores

ADD labels and legend

DISPLAY plot

FUNCTION visualize\_topics(lda\_model, dictionary, corpus, num\_topics, num\_words):

GENERATE word clouds for all topics

GENERATE bar charts for all topics

CREATE intertopic distance map using t-SNE

**DISPLAY** visualizations

#### MAIN PROCESS:

INITIALIZE empty lists for processed\_texts and raw\_texts

FOR EACH PDF file in pdf\_folder\_path:

text = extract\_text\_from\_pdf(file)

ADD text to raw texts

 $processed\_text = preprocess(text)$ 

ADD processed text to processed texts

FOR n FROM 2 TO 5:

COMPUTE and PRINT top 10 n-grams

CREATE dictionary from processed\_texts

CREATE corpus from processed\_texts and dictionary

best\_model, best\_num\_topics, best\_coherence, best\_perplexity = evaluate\_lda\_model(corpus, dictionary, processed\_texts)

CALL visualize\_topics with best\_model

PRINT best model's topics

SAVE best model to file

GENERATE and DISPLAY pyLDAvis visualization

### 5 LDA with Stemming

Algorithm 4 describes the pre-processing function for the LDA model with stemming. The remainder of the steps is equivalent to Algorithm 3. The respective code is available in LDA\_stemming\_final\_version.ipynb.

#### **Algorithm 4** Function preprocess(text) with Stemming

FUNCTION preprocess(text):

CONVERT text to lowercase

REMOVE numbers from text

REMOVE punctuation characters from text

TOKENIZE the lowercase text into individual words

INITIALIZE empty list for processed\_tokens

FOR EACH word in tokens:

IF word consists only of alphabetic characters AND word is not in stop\_words: stem = APPLY stemming to word ADD stem to processed\_tokens

RETURN processed\_tokens

#### **6** LDA with Translation

Algorithm 5 details the cross-language model in which the reports were translated to English before performing LDA. The respective code is available in the file translation\_lda\_final\_version.ipynb.

The functions <code>ensure\_json\_serializable</code>, <code>clean\_vis\_data\_in\_place</code> and <code>remove\_complex\_numbers</code> were created to solve the error TypeError: Object of type complex is not JSON serializable that was observed when trying to plot pyLDAvis.

#### Algorithm 5 Topic Modelling with LDA and Translation

SET seed to 23157003

#### INITIALIZE random number generators with seed

DOWNLOAD NLTK resources (punkt, stopwords, wordnet, omw-1.4)

DEFINE custom\_stopwords

COMBINE custom\_stopwords with English stopwords from NLTK

SET pdf\_directory to location of PDF files

SET translated\_texts\_file to "translated\_texts.pkl"

FUNCTION extract\_text\_from\_pdfs(pdf\_directory):

INITIALIZE empty list text\_data

FOR EACH PDF file in pdf\_directory:

TRY:

OPEN PDF file

FOR EACH page in PDF:

EXTRACT text from page

ADD text to text\_data

CATCH any exceptions:

PRINT error message

RETURN text data

FUNCTION translate\_text(text\_data, src\_lang="pt", tgt\_lang="en"):

LOAD pre-trained translation model and tokenizer

INITIALIZE empty list translated\_texts

FOR EACH text in text data:

TRY:

TOKENIZE and ENCODE text

GENERATE translation

DECODE translated tokens

ADD translated text to translated\_texts

CATCH any exceptions:

PRINT error message

ADD empty string to translated\_texts

RETURN translated\_texts

FUNCTION preprocess\_text(text):

CONVERT text to lowercase

REMOVE numbers

REMOVE punctuation

TOKENIZE text

REMOVE stopwords and short words

RETURN processed tokens

FUNCTION process\_texts(pickle\_file\_path):

LOAD texts from pickle file

INITIALIZE empty list processed texts

FOR EACH text in texts:

processed\_tokens = preprocess\_text(text)

ADD processed\_tokens to processed\_texts

RETURN processed\_texts

FUNCTION compute\_coherence\_and\_perplexity(dictionary, corpus, texts, start, limit, step):

INITIALIZE empty lists for coherence\_values, perplexity\_values, and model\_list

FOR num\_topics FROM start TO limit STEP step:

CREATE LdaModel with num topics

COMPUTE coherence score

COMPUTE perplexity

ADD scores and model to respective lists

RETURN model\_list, coherence\_values, perplexity\_values

FUNCTION evaluate\_lda\_model(corpus, dictionary, processed\_texts, start, limit, step):

CALL compute\_coherence\_and\_perplexity

NORMALIZE coherence and perplexity scores

COMPUTE composite scores

FIND best model based on highest composite score

PRINT best model details

PLOT evaluation metrics

RETURN best\_model, best\_num\_topics, best\_coherence, best\_perplexity

FUNCTION plot\_evaluation\_metrics(start, limit, step, coherence\_values, perplexity\_values, composite\_scores):

CREATE plot with three y-axes

PLOT coherence scores

PLOT perplexity scores

PLOT composite scores

ADD labels and legend

DISPLAY plot

FUNCTION print\_lda\_topics(lda\_model, num\_words):

FOR EACH topic in lda\_model:

PRINT topic words and weights

FUNCTION ensure ison serializable(data):

RECURSIVELY convert non-JSON-serializable data types to serializable types RETURN JSON-serializable data

FUNCTION clean\_vis\_data\_in\_place(vis\_data):

CLEAN 'Freq' column in vis\_data.token\_table using remove\_complex\_numbers

CLEAN 'Freq' column in vis\_data.topic\_info using remove\_complex\_numbers

CLEAN 'x' column in vis\_data.topic\_coordinates using remove\_complex\_numbers

CLEAN 'y' column in vis\_data.topic\_coordinates using remove\_complex\_numbers

FUNCTION remove\_complex\_numbers(data\_column):

INITIALIZE empty cleaned column

FOR EACH value IN data column:

IF value IS complex number:

ADD real part of value to cleaned\_column

ELSE:

ADD value to cleaned\_column

RETURN cleaned\_column

FUNCTION visualize\_lda\_model(lda\_model, corpus, dictionary, notebook):

PREPARE LDA visualization data

IF in notebook environment:

DISPLAY visualization in notebook

ELSE:

SAVE visualization as HTML file

OPEN HTML file in web browser

FUNCTION get\_top\_ngrams(processed\_texts, n, top\_k):

GENERATE n-grams from processed texts

COUNT frequency of n-grams

RETURN top\_k most common n-grams

FUNCTION save\_translated\_texts(texts, filename):

SAVE texts to pickle file

FUNCTION load\_translated\_texts(filename):

LOAD texts from pickle file

**RETURN** loaded texts

#### MAIN PROCESS:

IF translated texts file exists:

LOAD translated texts from file

ELSE:

 $text\_data = extract\_text\_from\_pdfs(pdf\_directory)$ 

translated\_texts = translate\_text(text\_data)

SAVE translated texts to file

texts = process\_texts(translated\_texts\_file)

CREATE dictionary from texts

CREATE corpus from texts and dictionary

best\_model, best\_num\_topics, best\_coherence, best\_perplexity = evaluate\_lda\_model(corpus, dictionary, texts)

PRINT details of best model

VISUALIZE best LDA model

FOR n FROM 2 TO 5:

COMPUTE and PRINT top 5 n-grams

## 7 Word2Vec + K-means

The model which combined word2vec and k-means was based on the tutorial given by Castillo (2018). The model is described in Algorithm 6 and its respective code can be found in file word2vec\_final\_version.ipynb.

#### Algorithm 6 Topic Modelling with Word2Vec and K-means

SET seed to 23157003

INITIALIZE random number generators with seed

**DOWNLOAD NLTK resources** 

SET pdf\_folder\_path to location of PDF files

DEFINE extra stop words

COMBINE extra\_stop\_words with Portuguese stopwords from NLTK

FUNCTION extract\_text\_from\_pdf(pdf\_path):

INITIALIZE empty text string

TRY:

OPEN pdf\_file

FOR EACH page in pdf\_file:

ADD page text to text string

CATCH any exceptions:

PRINT error message

RETURN text string

FUNCTION preprocess(text):

CONVERT text to lowercase

REMOVE numbers from text

REMOVE punctuation from text

TOKENIZE text

RETURN list of tokens not in stop\_words and longer than 2 characters

FUNCTION evaluate\_word2vec(model, words):

FOR EACH word in words:

TRY:

FIND most similar words using model

PRINT similar words and their scores

CATCH KeyError:

PRINT word not in vocabulary

FUNCTION vectorize(list\_of\_docs, model):

FOR EACH document in list of docs:

INITIALIZE zero vector

FOR EACH token in document:

IF token in model vocabulary:

ADD token vector to document vectors

IF document vectors not empty:

COMPUTE average of document vectors

ELSE:

USE zero vector

RETURN list of document vectors

FUNCTION mbkmeans\_clusters(X, k, mb, print\_silhouette\_values):

PERFORM MiniBatchKMeans clustering

COMPUTE silhouette score

IF print\_silhouette\_values is True:

COMPUTE and PRINT silhouette values for each cluster

RETURN clustering model, cluster labels, and average silhouette score

FUNCTION find\_best\_num\_topics(vectorized\_docs, min\_topics, max\_topics, step, mb):

INITIALIZE best\_k and best\_silhouette

FOR k FROM min\_topics TO max\_topics STEP step:

PERFORM clustering with k clusters

IF silhouette score is better than best\_silhouette:

UPDATE best k and best silhouette

PRINT best number of clusters and its silhouette score

PLOT silhouette scores vs number of clusters

RETURN best\_k

FUNCTION plot\_word\_embeddings(model, words):

FILTER words in model vocabulary

GET word vectors PERFORM t-SNE dimensionality reduction PLOT words in 2D space

FUNCTION plot\_clusters(vectorized\_docs, cluster\_labels):

PERFORM t-SNE dimensionality reduction on vectorized\_docs

PLOT documents in 2D space, colored by cluster labels

#### MAIN PROCESS:

INITIALIZE empty list for processed\_texts

FOR EACH PDF file in pdf\_folder\_path:

text = extract\_text\_from\_pdf(file)

 $processed\_text = preprocess(text)$ 

ADD processed\_text to processed\_texts

TRAIN Word2Vec model on processed\_texts

EVALUATE Word2Vec model with sample words

vectorized\_docs = vectorize(processed\_texts, word2vec\_model)

best\_num\_topics = find\_best\_num\_topics(vectorized\_docs, min\_topics=2, max\_topics=51, step=1, mb=100)

PERFORM clustering with best\_num\_topics

CREATE DataFrame with original texts, tokens, and cluster labels

PLOT word embeddings for sample words

**PLOT** clusters

PRINT top terms per cluster based on centroids

PRINT most frequent terms per cluster

ANALYZE and PRINT representative documents for a sample cluster

SAVE Word2Vec model

### 8 BERTopic

The BERTopic model is detailed in Algorithm 7. Its respective code can be found in the file named bertopic\_final\_version.ipynb.

#### Algorithm 7 Topic Modelling with BERTopic

SET seed to 23157003

INITIALIZE random number generators with seed

DOWNLOAD NLTK resources (stopwords and punkt)

SET pdf\_path to location of PDF files FUNCTION extract\_text\_from\_pdf(pdf\_file): INITIALIZE empty text string TRY: OPEN pdf\_file FOR EACH page in pdf\_file: ADD page text to text string CATCH any exceptions: PRINT error message RETURN text string FUNCTION preprocess text(text): CONVERT text to lowercase REMOVE numbers from text REMOVE punctuation from text TOKENIZE text into words REMOVE stopwords from tokens RETURN joined tokens as string FUNCTION compute coherence score(model, texts, topics, top n words): EXTRACT top words for each topic CREATE dictionary from texts COMPUTE coherence score using c\_v measure RETURN coherence score FUNCTION plot\_all\_topics(df, topic\_model): COMPUTE mean positions of each topic **GENERATE** colors for topics CREATE scatter plot of all topics ANNOTATE top 50 topics ADJUST text annotations to avoid overlap DISPLAY plot FUNCTION plot\_top\_10\_topics(df, topic\_model): IDENTIFY 10 most frequent topics FILTER DataFrame for top 10 topics COMPUTE mean positions of each topic **GENERATE** colors for topics CREATE scatter plot of top 10 topics ANNOTATE top 10 topics ADJUST text annotations to avoid overlap DISPLAY plot FUNCTION visualize hierarchy custom(topic model, topics, top n topics, width, height):

IF topics not provided:

GET all topics from model

IF top\_n\_topics specified:

CREATE labels for each topic

LIMIT topics to top\_n\_topics

**EXTRACT** topic embeddings

PERFORM hierarchical clustering

CREATE labels for topics

GENERATE dendrogram

#### DISPLAY plot

FUNCTION visualize\_topics\_barchart(topic\_model, top\_n\_topics, n\_words, width, height):

GET topic information from model

SORT topics by size (excluding outliers)

SELECT top N topics

CREATE labels for each topic

GENERATE horizontal bar chart

CUSTOMIZE plot appearance

ADD count labels to bars

DISPLAY plot

#### INITIALIZE empty list for texts

FOR EACH pdf\_file in pdf\_path:

extracted\_text = extract\_text\_from\_pdf(pdf\_file)

IF extracted\_text is not empty:

preprocessed\_text = preprocess\_text(extracted\_text)

ADD preprocessed text to texts list

DEFINE extra\_stop\_words

COMBINE extra\_stop\_words with Portuguese stopwords

CREATE CountVectorizer with combined stopwords

CREATE UMAP model for dimension reduction

INITIALIZE empty list for coherence\_scores

SET topic\_range to range from 2 to 50, step 2

FOR EACH num\_topics in topic\_range:

CREATE BERTopic model with num\_topics

FIT model to texts

COMPUTE coherence score

ADD coherence score to coherence scores list

PLOT coherence scores against number of topics

SET optimal\_topics to number of topics with highest coherence score

CREATE final BERTopic model with optimal\_topics

FIT final model to texts

SAVE final model

REDUCE dimensionality of embeddings to 2D using UMAP

ENSURE reduced\_embeddings is 2D

CREATE DataFrame with reduced embeddings, topics, and text lengths

**GENERATE** visualizations:

CALL plot\_all\_topics()

CALL plot\_top\_10\_topics()

CALL visualize\_hierarchy\_custom()

CALL visualize\_topics\_barchart()

# 9 Stop Words

The list of stop words used to customise the NLTK file is provided in Table 2. The model which performed LDA to the translated documents used the NLTK file for the English language. For this reason, the stop words list was customised using the translated version of the words.

**Table 2: Extra Stop Words** 

Portuguese Words	English Words	Portuguese Words	English Words
acidentes	accidents	empresa	enterprise
acerca	about	errôneas	erroneous
acidente	accident	exclusivamente	exclusively
acordo	agreement	fabricante	manufacturer
aeronáutica	aeronautics	fatos	facts
aeronáuticas	aeronautical	figura	figure
aeronáutico	aeronautic	final	final
aeronáuticos	aeronautical	fins	purposes
aeronave	aircraft	futuros	future
aeronaves	aircraft	geral	general
agência	agency	glossário	glossary
ambos	both	havia	had
anac	anac	hipóteses	hipothesis
and	and	incidente	incident
anexo	attachment	induzir	induce
após	after	informações	information
aspectos	aspects	internacional	international
aviação	aviation	interpretações	Interpretations
avião	airplane	investigação	research
brasileiro	brazilian	investigações	investigations
cenipa	cenipa	junto	together
centro	center	matrícula	registration
civil	civil	modelo	model
conclusão	conclusion	momento	moment
conduzidas	conducted	nacional	national
conforme	according	nada	nothing
contribuintes	taxpayers	neste	this
contribuiu	contributed	n°	no
cuja	whose	nota	note
culpa	blame	nsca	nsca
dados	data	objetivo	objective
desconhecido	unknown	ocorrência	occurrence
deste	this	of	of
determinar	determine	organização	organization

Portuguese Words	English Words		Portuguese Words	English Words
organizacionais	organizational		sipaer	sipaer
outras	other		sobre	upon
outro	other		suma	all
outros	other		tampouco	either
piloto	pilot		ter	have
pode	can		termos	terms
poderá	can		the	the
poderá	can		tipo	kind
possível	possible		to	to
possuia	had		totais	total
prevenção	prevention		trazer	bring
propósito	purpose		últimas	last
qualquer	any		últimos	last
realizar	accomplish		uso	use
registros	records		utc	utc
relatar	tell		utilização	use
relatório	report		V00	flight
segurança	safety		vôo	flight
seripa	seripa	'		
simplificado	simplified			

# References

Castillo, D. (2018) *How to Cluster Documents Using Word2Vec and K-Means*. Available at: https://dylancastillo.co/posts/nlp-snippets-cluster-documents-using-word2vec.html.