

# Intelligent Classification of Legislative Proposals

Luca Peres Quinta da Guarda   Ronie Paulucio Porfirio

**Universidade de Brasília (UnB)**

*Programa de Pós-Graduação em Computação Aplicada*

**Course Mineração de Dados Massivos**

*Prof. Dr. Marcelo Ladeira and Pr. MSc. Gustavo Cordeiro*

# Contents

- 1 Introduction
- 2 The Problem
- 3 Objective
- 4 Literature review
- 5 Methodology
- 6 Experiments Conducted
- 7 Conclusions
- 8 Trabalhos Futuros

# Introduction

## Business Understanding - Legislative Proposal's Definition and Types

### Legislative Proposal

- A **Legislative Proposal** is any matter subject to deliberation by the Legislative Chamber (CLDF Internal Regulations, art. 129).

### Legislative Proposals can be of the following types:

- Proposta de emenda à Lei Orgânica;
- Projeto de lei complementar;
- Projeto de lei;
- Projeto de decreto legislativo;
- Projeto de resolução;
- Indicação;
- Moção;
- Requerimento;
- Emenda;
- Recursos;

# Introduction

## Business Understanding - Legislative Proposal's Themes

Legislative Proposals are classified into one or more themes:

- Agricultura
- Assistência Social
- Assunto Fundiário e Ordenamento Territorial
- Assunto Social
- Cidadania
- Ciência e Tecnologia
- Combate à Corrupção
- Comunicação
- Comércio e Serviços
- Cultura
- Defesa do Consumidor
- Desenvolvimento Econômico
- Desporto e Lazer
- Direitos Humanos
- Economia
- Educação
- Energia
- Fiscalização e Governança
- Habitação
- Incentivos Fiscais e Concessões Públicas
- Indústria
- Meio Ambiente
- Não se aplica
- Outro
- Previdência Social
- Relações Exteriores
- Saneamento
- Saúde
- Segurança
- Servidor Público
- Trabalho
- Transporte e Mobilidade Urbana
- Turismo
- Urbanismo

**Total: 34 themes**

# Introduction

## Business Understanding - Why Themes?

### Thematic Classification Benefits

- Efficient classification of legislative proposals is crucial to **streamline their analysis** and processing within the legislative process helping to **determine which committees a proposal should go through**.
- By categorizing legislative proposals into relevant themes, lawmakers can streamline their analysis, **allocate resources efficiently**, and **make informed decisions**.
- This process **enhances transparency** and facilitates a more **organized legislative workflow**.
- Thematic classification plays an important role in maintaining **accurate information retrieval** and **ensuring effective legislative management**.

# The Problem

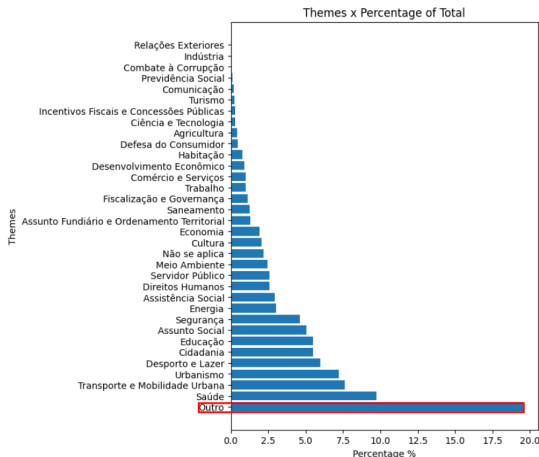
## Understanding the problem

### Theme “others” is growing bigger

- Usually, the author of the proposal is responsible for classifying it into one or more themes.
- Unfortunately, due to various factors such as ambiguous topics, outdated categories, multidisciplinary nature, **many propositions end up being classified under the generic label of “others”**.

# The Problem

## Understanding the problem



The chart shows that the number of proposals classified under the theme “others” represents approximately 20% of the total.

# The Problem

## Problem definition

### The problem is inadequate proposal classification

Inadequate classification hinders efficient tracking, analysis, and transparency of legislative activities, making it difficult for both society and lawmakers to understand and oversee the legislative process effectively.





# Objectives

## Primary objective

### Primary objective

- The primary objective of this study is to **compare different machine learning models** to determine which model is **most effective in suggesting more appropriate theme categories for legislative proposals**.
- The goal is to find **more suitable themes that better match the content of these proposals**.



# Objectives

## Disclaimer

### Attention!

- This study **does not aim to automate the classification process, replace human classification, or compare human classification with machine learning classification.**
- Instead, **the focus is on enhancing the existing categorization process by identifying the best model for suggesting themes for proposals currently classified under the generic label "others."**

# Literature review



Reimers, Nils and Gurevych, Iryna (2019)

Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks

*Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing Association for Computational Linguistics*



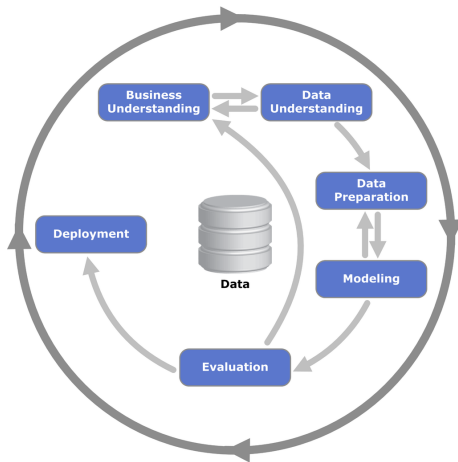
J. Andrade Junior, J. Cardoso-Silva, and L. Bezerra (2021)

Comparing Contextual Embeddings for Semantic Textual Similarity in Portuguese

*Anais da X Brazilian Conference on Intelligent Systems*

# Methodology

## CRISP-DM's Methodology



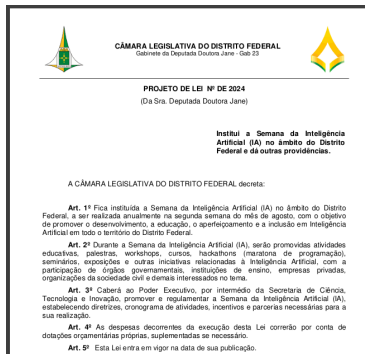
We utilized the CRISP-DM methodology. Its stages will be explained in the following slides.

# Methodology

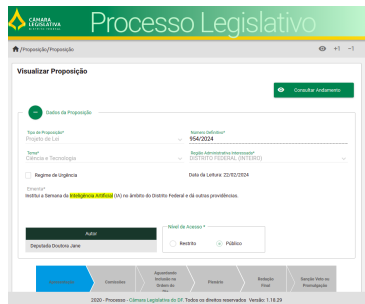
## 1 - Business Understanding

### 1 - Business Understanding

We have analyzed legislative documents to align our data mining objectives with legislative classification needs.



(a) Proposal example



(b) Proposal's data

# Methodology

## 2 - Data Understanding

### 2 - Data Understanding

The dataset contains 22,267 summaries extracted from the Electronic Legislative Process (PLE) system, covering the period from 2021 to May 2024. Each summary is accompanied by its respective thematic classification.



### Data Preparation

- ① **Preprocessing:** We discard data classified under “outro” and “não se aplica” themes. Then, we perform tokenization, normalization, stopwords removal, and lemmatization processes.
- ② **Vectorization:**
  - **Multilingual sentence embedding model** based on the MiniLM architecture, a lightweight and efficient BERT variant with 12 transformer layers, to produce high-quality embeddings that capture the semantic meaning of the text.
  - **TF-IDF (Term Frequency-Inverse Document Frequency) model** involves converting text into numerical vectors based on the frequency of terms within a document and across a collection of documents. This approach captures the importance of a term in a document relative to the entire dataset.

### 4 - Modeling

- 1 **DummyClassifier:** A **baseline model** that makes predictions using simple rules and establishes a baseline to **compare the performance of more complex models**.
- 2 **Support Vector Machine (SVM):** A powerful model that **finds the hyperplane that best separates the classes in the feature space**. It is effective for high-dimensional spaces and when the number of dimensions exceeds the number of samples.
- 3 **Logistic Regression:** A linear model that **estimates the probability of a binary outcome based on input features**. It is simple and interpretable, good for linearly separable data and understanding feature importance.



### 4 - Modeling (cont...)

- 4 **XGBoost (Extreme Gradient Boosting):** An optimized gradient boosting algorithm that **builds an ensemble of weak learners (typically decision trees) to improve model performance**. Known for high performance, speed, and scalability.
- 5 **Random Forest:** An **ensemble model that constructs multiple decision trees and aggregates their predictions**. Robust against overfitting, good for handling large datasets with higher dimensionality.
- 6 **K-Nearest Neighbors (KNN):** A **non-parametric model that classifies a data point based on the majority class of its k nearest neighbors**. Simple and intuitive, effective for small datasets with low noise.

### Why Choose These Models?

- 1 **Baseline Comparison:** DummyClassifier provides a benchmark to gauge the performance of more sophisticated models.
- 2 **Linear and Non-Linear Data:** Logistic Regression and SVM cover linear relationships, while Random Forest, XGBoost, and KNN handle non-linear data.
- 3 **Model Performance:** XGBoost and Random Forest are chosen for their strong predictive performance and ability to handle complex datasets.
- 4 **Interpretability:** Logistic Regression is valued for its simplicity and ease of interpretation.
- 5 **Versatility:** SVM, Random Forest, and XGBoost offer versatility across different types of data and problems.
- 6 **Scalability:** XGBoost is particularly chosen for its scalability and efficiency in handling large datasets.

### 5 - Evaluation

Our evaluation metrics include **Accuracy, Precision, Recall** and **F1 score**.

- Accuracy is **straightforward and easy to understand**.
- Precision and recall are more informative when dealing with **imbalanced classes** because they provide insights into the performance of the minority class, which accuracy might overlook.
- F1 score gives a **balance between precision and recall as a single metric** to summarize the model performance.

### 6 - Deployment

- The results might be used to create an interface in the PLE system to suggest themes that best fit new proposals.

# Experiments Conducted

nao sei

	Model	fit.time	score.time	test.f1.weighted	test.balanced.accuracy	test.precision.weighted	test.recall.weighted
0	DummyClassifier	0.006086	<b>0.016813</b>	0.025819	0.031250	0.014462	0.120258
1	KNeighborsClassifier	<b>0.005841</b>	33.463293	0.323486	0.182068	0.447307	0.299509
2	<b>LogisticRegression</b>	5.160993	0.024248	0.467718	<b>0.377732</b>	0.506144	0.458011
3	RandomForestClassifier	373.229191	1.552665	0.435226	0.273611	0.444914	0.449233
4	SVM	52.730176	14.002061	<b>0.472803</b>	0.319393	<b>0.512542</b>	0.456292
5	XGBClassifier	128.124318	1.915238	0.449333	0.277803	0.449830	<b>0.458809</b>

Table: nao sei

Não sei

- Explicar tabela

# Experiments Conducted

## Embedding - Cross Validation

	fit_time	score_time	test_f1_weighted	train_f1_weighted	test_balanced_accuracy	train_balanced_accuracy	test_precision_weighted	train_precision_weighted	test_recall_weighted	train_recall_weighted
Model										
DummyClassifier	0.028611	0.013967	0.025819	0.025819	0.031250	0.031250	0.014462	0.014462	0.120258	0.120258
KNeighborsClassifier	0.040055	1.049414	0.307679	0.516615	0.179927	0.370583	0.323286	0.547667	0.311971	0.520135
LogisticRegression	3.597823	0.025262	0.252122	0.282075	0.217896	0.387424	0.356856	0.397779	0.234561	0.270780
RandomForestClassifier	225.171689	0.825105	0.299710	<b>0.807638</b>	0.163765	<b>0.865793</b>	0.328777	<b>0.835075</b>	0.327624	0.801934
SVM	42.684947	33.030448	0.295933	0.346580	<b>0.241816</b>	0.495388	<b>0.376351</b>	0.435997	0.276796	0.337201
XGBClassifier	253.781450	0.560534	<b>0.325231</b>	0.800150	0.178569	0.726839	0.330270	0.800984	<b>0.341866</b>	<b>0.802977</b>

Table: Embedding - Cross Validation

## Embedding

- Explicar tabela

# Experiments Conducted

## TfidfVectorizer - Cross Validation

	fit_time	score_time	test_f1_weighted	train_f1_weighted	test_balanced_accuracy	train_balanced_accuracy	test_precision_weighted	train_precision_weighted	test_recall_weighted	train_recall_weighted
Model										
DummyClassifier	0.005016	0.013737	0.025819	0.025819	0.031250	0.031250	0.014462	0.014462	0.120258	0.120258
KNeighborsClassifier	0.004096	28.906825	0.323249	0.541327	0.181929	0.396655	0.446675	0.629905	0.299141	0.527901
LogisticRegression	<b>3.725738</b>	<b>0.018837</b>	<b>0.467718</b>	<b>0.584556</b>	<b>0.377732</b>	<b>0.727172</b>	<b>0.506144</b>	<b>0.629863</b>	<b>0.458011</b>	<b>0.585697</b>
RandomForestClassifier	267.923079	1.148535	0.435199	<b>0.806360</b>	0.272152	<b>0.864809</b>	0.445313	<b>0.834005</b>	0.449355	<b>0.801044</b>
SVM	44.651407	12.235956	<b>0.472803</b>	0.676972	0.319393	0.803307	<b>0.512542</b>	0.712173	0.456292	0.673389
XGBClassifier	102.740959	1.606697	0.449333	0.766095	0.277803	0.768935	0.449830	0.767603	<b>0.458809</b>	0.769460

Table: TfidfVectorizer - Cross Validation

## TfidfVectorizer

- Explicar tabela

# Experiments Conducted

## TfidfVectorizer and MultiOutputClassifier - Cross Validation

	fit_time	score_time	test_f1_weighted	train_f1_weighted	test_balanced_accuracy	train_balanced_accuracy	test_precision_weighted	train_precision_weighted	test_recall_weighted	train_recall_weighted
Model										
DummyClassifier	0.004188	0.011938	0.025819	0.025819	0.031250	0.031250	0.014462	0.014462	0.120258	0.120258
KNeighborsClassifier	0.005413	29.213731	0.323249	0.541327	0.181929	0.396655	0.446675	0.629905	0.299141	0.527901
LogisticRegression	3.685860	0.019390	0.467718	0.584556	<b>0.377732</b>	0.727172	0.506144	0.629863	0.458011	0.585697
RandomForestClassifier	268.849362	1.154890	0.434992	<b>0.806396</b>	0.272071	<b>0.864844</b>	0.444953	<b>0.834183</b>	0.449048	<b>0.800982</b>
SVM	44.404848	12.142895	<b>0.472803</b>	0.676972	0.319393	0.803307	<b>0.512542</b>	0.712173	0.456292	0.673389
XGBClassifier	102.840638	1.657798	0.449333	0.766095	0.277803	0.768935	0.449830	0.767603	<b>0.458809</b>	0.769460

Table: TfidfVectorizer and MultiOutputClassifier - Cross Validation

## TfidfVectorizer and MultiOutputClassifier

- Explicar tabela



# Conclusions

## Conclusions

- Quais?

# Future Work

## Future Work

- Quais?