



ACH2118

Introdução ao Processamento de Língua Natural

---

**Relatório de Implementação do Classificador de Clareza para Respostas na  
Plataforma eSIC**

---

**Aluno**

Felipe Mateos Castro de Souza - 11796909

Nilton Tadashi Enta - 12730911

**Docente**

Ivandr  Paraboni

Escola de Artes, Ci ncias e Humanidades - EACH

13 de novembro de 2023

# Sumário

<b>1</b>	<b>Introdução</b>	<b>2</b>
<b>2</b>	<b>Análise Exploratória</b>	<b>2</b>
<b>3</b>	<b>Pré-processamento</b>	<b>3</b>
3.1	Tratamento dos dados . . . . .	3
3.2	Separação do Conjunto de Dados . . . . .	3
3.3	Técnicas de Vetorização de Texto . . . . .	4
3.4	Técnicas de Word Embeddings . . . . .	4
<b>4</b>	<b>Treinamento de modelos</b>	<b>4</b>
4.1	Teste de Modelos de classificadores . . . . .	4
4.2	Modelagem do Algoritmo . . . . .	4
<b>5</b>	<b>Avaliação dos Modelos</b>	<b>5</b>
5.1	Conjunto de treinamento . . . . .	5
5.2	Conjunto de teste . . . . .	5
<b>6</b>	<b>Reprodução dos resultados</b>	<b>5</b>

# 1 Introdução

Este relatório apresenta os resultados da implementação de um classificador de nível de clareza para respostas publicadas na plataforma eSIC. O objetivo central do projeto foi desenvolver um modelo capaz de categorizar as respostas em três níveis de clareza, identificados como 'c1', 'c234' e 'c5'. O trabalho foi conduzido em duas partes distintas: a primeira envolveu o desenvolvimento do classificador e a apresentação dos resultados médios de acurácia utilizando validação cruzada de 10 partições sobre o conjunto de dados de treinamento. A segunda parte consistiu na geração de previsões para o conjunto de teste.

Ao longo deste processo, realizamos uma análise exploratória dos dados, aplicamos técnicas de pré-processamento, vetorização de texto e implementamos modelos de aprendizado de máquina clássicos e neurais. O relatório detalha cada etapa do desenvolvimento, desde a importação e análise dos dados até a avaliação final dos modelos no conjunto de teste.

Os algoritmos testados incluíram Naive Bayes, Support Vector Machine, Random Forest e LSTM, sendo que a otimização dos modelos foi realizada através do grid search para SVM e Random Forest. As métricas de avaliação, como acurácia média na validação cruzada e acurácia no conjunto de teste, foram registradas para cada algoritmo, proporcionando uma visão abrangente do desempenho de cada modelo.

Os resultados obtidos neste relatório são essenciais para a escolha do modelo mais adequado para a tarefa de classificação de respostas na plataforma eSIC, contribuindo para a eficácia na análise de clareza dessas respostas.

## 2 Análise Exploratória

Realizamos uma análise preliminar do conjunto de dados utilizando as bibliotecas pandas e numpy.

Verificamos as informações gerais da base como presença de valores nulos, tipo de cada feature e quantidade de linhas. O dataframe possui 6 mil linhas e nenhuma vazia, todas do tipo 'object'.

Em seguida agrupamos os dados por classe e partindo da função 'describe' verificamos a presença de dados duplicados. Cada classe possui um conjunto de 2 mil textos sendo únicos:

c1 - 1.839

c234 - 1.929

c5 - 1911

E analisamos também as palavras que se repetiam em cada classe. Mapeamos as dez maiores classes e em comum ambas tinham como as três palavras mais presente em seus textos: 'informação', 'por' e 'para'.

## 3 Pré-processamento

### 3.1 Tratamento dos dados

No âmbito do tratamento de dados textuais, a análise precisa e eficiente requer uma série de técnicas avançadas. Abaixo descrevemos o processo que utilizamos, abordando desde a tokenização até a codificação de rótulos.

a. Tokenização: A primeira etapa envolve a tokenização, um processo fundamental que divide o texto em unidades semânticas chamadas tokens. Cada palavra ou expressão é isolada, facilitando análises subsequentes e fornecendo uma visão granular do conteúdo textual. Optamos por utilizar a biblioteca NLTK para esta tarefa.

b. Lemmatização e remoção de pontuações: A lematização entra em cena para simplificar a análise ao reduzir palavras flexionadas às suas formas base. Essa técnica garante que diferentes formas de uma palavra sejam tratadas como uma única entidade, facilitando a compreensão e reduzindo a complexidade do conjunto de dados. Já para esta tarefa, optamos por usar a biblioteca SpaCy, onde fizemos uso do modelo 'pt\_core\_news\_sm', vale ressaltar que juntamente desse processo, fizemos a remoção de pontuações, usando o tagger da biblioteca SpaCy.

c. POS Tagging (Etiquetagem de Partes da Fala): A etiquetagem de partes da fala adiciona camadas de informações valiosas, atribuindo etiquetas a cada palavra de acordo com sua função gramatical. Essa abordagem aprimora a compreensão contextual, permitindo uma análise mais profunda das relações sintáticas e semânticas entre as palavras.

d. 2-Grams e 3-Grams: A implementação de bigramas (2-grams) e trigramas (3-grams) expande a análise para além das palavras isoladas, considerando sequências de duas e três palavras consecutivas. Isso captura relações mais complexas e contextuais no texto, fornecendo insights mais ricos sobre a estrutura e o significado. Para essa tarefa optamos por utilizar a mesma biblioteca que no tokenizer, isto é, a NLTK.

e. Label Encoding (Codificação de Rótulos): A etapa final do processo consiste na codificação de rótulos, transformando informações textuais em representações numéricas. Essa técnica é crucial para a utilização eficiente de algoritmos de aprendizado de máquina, permitindo a aplicação de modelos preditivos a dados previamente tratados.

Para sermos mais eficientes, optamos por armazenar os tratamentos feitos em arquivos do tipo CSV, sendo o mais relevante deles o "esic2023\_cleaned.csv".

### 3.2 Separação do Conjunto de Dados

Com a finalidade de aumentar nossa confiança nos resultados das métricas que seriam obtidas após o treinamento dos modelos, optamos por dividir o conjunto de dados em duas partes, um conjunto de treino, correspondente a 70% dos dados, e um conjunto de teste, correspondente a 30% dos dados.

Dessa forma, podemos realizar uma média de uma validação cruzada no conjunto de treino

e depois comparar seus resultados com a métrica de acurácia no conjunto de teste, para aferir se houve *overfitting*.

### 3.3 Técnicas de Vetorização de Texto

Uma abordagem eficaz para essa tarefa é a utilização do método TF-IDF (Term Frequency-Inverse Document Frequency), e para implementar essa técnica de maneira robusta e eficiente, a biblioteca Scikit-Learn oferece ferramentas poderosas.

O Scikit-Learn fornece uma implementação flexível e fácil de usar do TF-IDF, permitindo que os usuários ajustem parâmetros conforme necessário e integrem facilmente essa etapa em seus pipelines de PLN e aprendizado de máquina. Desse modo optamos por utilizá-la em nosso projeto, com o hiperparâmetro 'max\_features=5000' na coluna 'lemma' dos conjuntos de treino e teste.

Optamos por armazenar as matrizes resultantes em arquivos NPZ e NPY, para os atributos e rótulos respectivamente, nos poupando, então, de ter que gerá-los novamente toda vez que reiniciássemos o kernel do Jupyter.

### 3.4 Técnicas de Word Embeddings

Uma abordagem inovadora para essa tarefa é o uso do modelo BERTimbau, uma variação do BERT (Bidirectional Encoder Representations from Transformers), presente na renomada biblioteca transformers. Optamos por utilizá-lo na seguinte configuração: return\_tensors="pt", truncation=True, padding=True, max\_length=512, add\_special\_tokens = True. Pois foi a que conferiu maior resultado dentre as abordagens usando embeddings,

## 4 Treinamento de modelos

### 4.1 Teste de Modelos de classificadores

Sob o data frame disponibilizado foram aplicados diversos classificadores, não neurais e neurais. Dentre os não neurais foi testado a Naive Bayes, Regressão Logística, Random Forest e SVM. Para os modelos neurais foi escolhido a LSTM e o BERT, ambos com uma rede neural como classificador.

As redes neurais demandam um tempo maior de processamento e apresentaram uma acurácia similar ou próxima das não neurais, porém com uma variação maior entre resultados. Sendo assim, escolhemos seguir com os algoritmos não neurais.

### 4.2 Modelagem do Algoritmo

Realizamos a otimização utilizando grid search do sklearn nos modelos de SVM e RF, escolhidos por ter a acurácia maior dentre os modelos não neurais. Após alguns testes e seleção de parâmetros o modelo que apresentou o melhor desempenho foi o SVM.

## 5 Avaliação dos Modelos

Após termos realizado a otimização de busca de hiperparâmetros usando o grid search da biblioteca sklearn nos algoritmos mais promissores, chegamos à conclusão que o algoritmo de SVM (Support Vector Machine) se saía melhor que os demais (Naive Bayes, Random Forest e LSTM), principalmente usando a vetorização produzida pelo TFIDF.

Desse modo, nosso *ensemble* final ficou com as seguintes características:

Vetorização: TFIDF dos textos lematizados e sem pontuação, com 'max\_features'=5000 Algoritmo: SVM com 'C'=2.0 e o 'kernel'=rbf

### 5.1 Conjunto de treinamento

Dentro do conjunto de treinamento, que corresponde a 70% dos do conjunto de dados, fizemos uma validação cruzada de 10 *folds*. Após a obtenção da média dos valores dessa métrica, constatamos que o modelo obteve uma acurácia média de 56.55%.

### 5.2 Conjunto de teste

Já no conjunto de teste, que corresponde a 30% dos dados, realizamos um teste de acurácia para aferir se o modelo permanecia com resultados consistentes. Obtivemos o resultado de 58.17% acurácia. Essa diferença aparenta ser pequena demais, e portanto, nos diz que o modelo não apresentou 'overfitting'.

## 6 Reprodução dos resultados

Basta seguir as instruções presentes no arquivo main.ipynb