

Manual de Utilizador - Classificador de Texto

Instituto Politécnico de Setúbal - ESTS

Metrado Engenharia Informática - Extração Automática de Informação

António Carlos Ferreira Pinto - 201801432

Diogo Costa

Guilherme Malhado



**POLITECNICO
SETUBAL**

Índice

1. [Introdução](#)
2. [Instalação](#) 2.1 [Requisitos](#) 2.2 [Correr a aplicação pela primeira vez:](#)
3. [Guia de Utilização](#) 3.1 [Estrutura da aplicação web:](#) 3.1.1 [Selects - Selecionar dados do nosso dataset](#)
3.1.2 [Process](#) 3.1.3 [Class Identifiers](#) 3.1.4 [Selects - Selecionar dados do nosso dataset](#)
4. [Anexos](#) 4.1 [Anexo I](#)

Introdução

Uma aplicação web simples, maioritariamente de página unica. É um classificador de textos, mais em específico reviews de comida, tendo duas classes para as reviews, reviews Positivas e reviews Negativas. Feito usando Node.js com Express.js, e uma base de dados em mysql

Instalação

Requisitos:

- Nodejs versão 18+

Correr a aplicação pela primeira vez:

1. Antes de correr o programa localmente criar um ficheiro .env usando o sample.env como base.
2. Num terminal na diretoria da aplicação correr:

```
npm i
```

3. No mesmo terminal correr para inicializar a aplicação em localhost, correr:

```
npm start
```

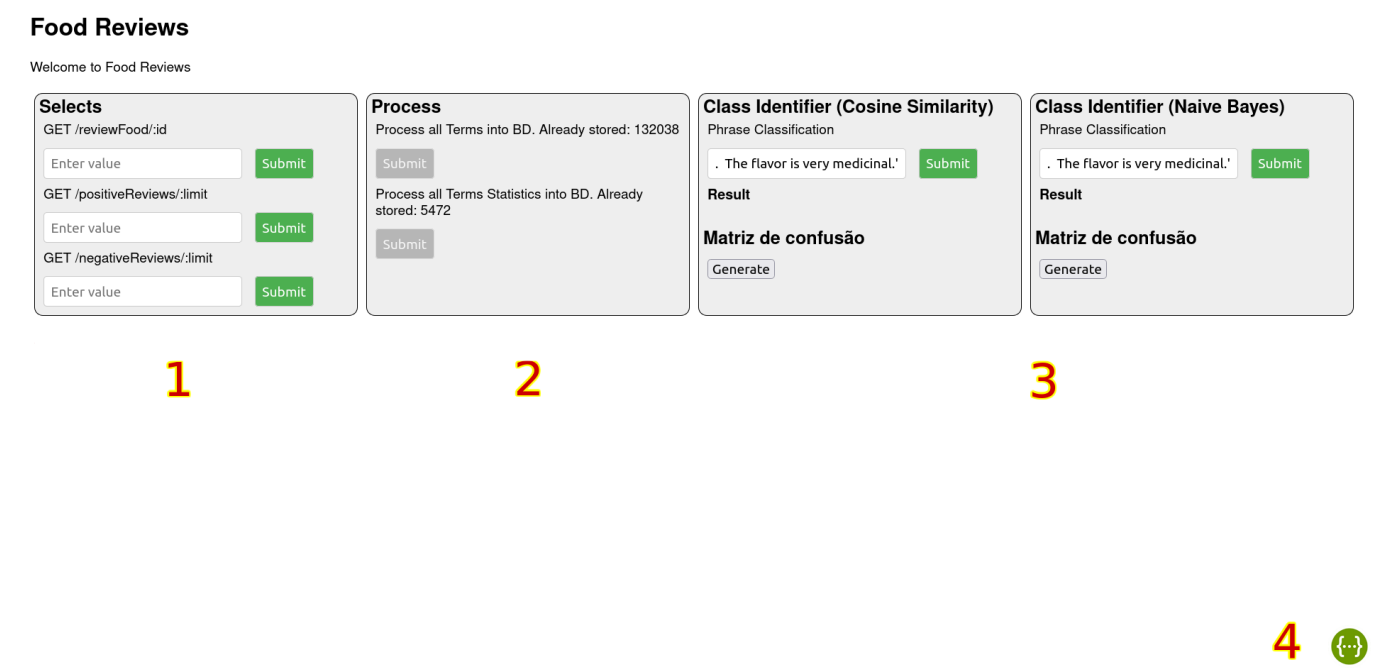
Guia de Utilização

Para inicializar inicializar o a aplicação correr:

```
npm start
```

Estrutura da aplicação web:

A aplicação funciona maioritariamente numa única página Composta por um titulo e 3 regiões principais distintas (enumeradas na image seguinte)



1. Selects - Selecionar dados do nosso dataset

Selects

GET /reviewFood/:id

GET /positiveReviews/:limit

GET /negativeReviews/:limit

Neste componente o utilizador pode facilmente fazer queries à base de dados para ver diferentes entradas do dataset utilizado para o nosso modelo.

Existem 3 alneas nesta região:

- GET /review/:id -> Endpoint que retorna a review com o id introduzido na textbox;
- GET /positiveReviews/:limit -> Endpoint que retorna X reviews classificadas como positivas (Score de 4 ou 5), X sendo o valor introduzido na textbox;
- GET /negativeReviews/:limit -> Endpoint semelhante ao anterior, mas para reviews negativas (Score de 1 ou 2).

2. Process

Esta região da aplicação é desaconcelhável, como tal esta funcionalidade está desativada, uma vez que quando qualquer um destes processos é inicializado irá demorar uma quantidade de tempo considerável. Para além de que não precisam de ser executados com regularidade.

Process

Process all Terms into BD. Already stored: 132038

Submit

Process all Terms Statistics into BD. Already stored: 5472

Submit

Nesta região temos dois botões:

- O primeiro para processar todos os dados contidos no nosso training set, de forma a obter os termos para os nossos classificadores;
- O segundo para depois agregar e processar as componentes desses termos.

Podemos ver nesta região a quantidade de entradas criadas na nossa tabela de termos, sendo cada entrada a presença de um termo num documento; podemos também ver o número de entradas na nossa tabela de estatísticas dos termos onde estão os nossos melhores termos agregados.

3. Class Identifiers

Temos duas secções semelhantes nesta região, funcionam de forma idêntica apenas aplicam algoritmos de classificação diferentes. O nome do algoritmo respetivo encontra-se no título da região, sendo estes, Similaridade de Cosseno e Naive Bayes.

Na "caixa" de cada secção temos dois elementos o classificador de frases e o gerador da nossa matriz de confusão.

Class Identifier (Cosine Similarity)

Phrase Classification

Result

Matriz de confusão

Class Identifier (Naive Bayes)

Phrase Classification

Result

Matriz de confusão

Secção 1 - Testar Frases/Reviews originais

Nesta secção o utilizador pode testar uma review de comida original para ver se é identificada com positiva ou negativa. É só escrever uma frase na textbox e clicar no botão de submit ao lado. O resultado será a resposta se é classificado como Positivo ou Negativo e os valores ponderados (Similarity values ou Product of Probabilities) respetivos a cada Classe (por ordem Positiva e Negativa).

Class Identifier (Cosine Similarity)

Phrase Classification

Result

Negative

Similarity values: 0.4617007324613109,
0.5123716578027154

Matriz de confusão

Class Identifier (Naive Bayes)

Phrase Classification

Result

Negative

Product of the probabilities:
6.893869280121e+43, 4.584144640657008e+44

Matriz de confusão

Secção 2 - Matriz de confusão - Testar o classificador

Nesta secção inicialmente temos só um botão para gerar os nossos testes, este quando primido gera a Matriz de confusão e as Métricas Prec, Rec e F1 (Ver anexo I para explicação de o que significa cada um) **A primeira vez que este processo é efetuado poderá demorar alguns minutos a calcular.**

Class Identifier (Cosine Similarity)

Phrase Classification

. The flavor is very medicinal.'

Submit

Result

Negative

Similarity values: 0.4617007324613109, 0.5123716578027154

Matriz de confusão

	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP) 102	False Negative (FN) 98
Actual Negative	False Positive (FP) 94	True Negative (TN) 106

Metricas

Prec:
0.5204081632653061

Rec:
0.51

F1:
0.5151515151515151

Class Identifier (Naive Bayes)

Phrase Classification

. The flavor is very medicinal.'

Submit

Result

Negative

Product of the probabilities:
6.893869280121e+43, 4.584144640657008e+44

Matriz de confusão

	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP) 153	False Negative (FN) 47
Actual Negative	False Positive (FP) 38	True Negative (TN) 162

Metricas

Prec:
0.8010471204188482

Rec:
0.765

F1:
0.7826086956521738

4. Swagger

[Docs](#)

Anexos

Anexo I

Matriz de confusão

Metricas

Prec - Precision

Precision é uma métrica que indica o quão preciso o modelo está de acordo com os que foram previstos positivamente, quantos deles é que são de facto positivos. É uma boa medida para determinar quando os custos de Falso Positivo são altos. [$\text{Precision} = \frac{\text{Positivos Verdadeiros}}{\text{Positivos Verdadeiros} + \text{Falsos Positivos}}$]

Rec - Recall

Recall calcula quantos dos Positivos Verdadeiros são capturados quando estimamos como Positivos Verdadeiros. Neste caso o Recall é uma métrica que terá grande importância quando há um alto custo associado aos Falsos Negativos. $\text{Recall} = \frac{\text{Positivos Verdadeiros}}{\text{Positivos Verdadeiros} + \text{Falsos Negativos}}$

$$[\text{Recall} = \frac{\text{Positivos Verdadeiros}}{\text{Positivos Verdadeiros} + \text{Falsos Negativos}}]$$

F1 - Score ou F-Measure

É utilizada quando se pretende encontrar um equilíbrio entre o Precision e o Recall.

$$[\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}]$$