

**FUNDAÇÃO GETULIO VARGAS**  
**MESTRADO EM MATEMÁTICA APLICADA**

**ALYNE PEREIRA DE OLIVEIRA**

**TREINAMENTO DE UM CLASSIFICADOR PARA RECONHECIMENTO  
DE CLASSES DE NOTÍCIAS**

**RIO DE JANEIRO**

**2022**

## **Sumário**

<b>1. INTRODUÇÃO</b>	<b>3</b>
<b>2. METODOLOGIA</b>	<b>5</b>
<b>3. RESULTADOS E DISCUSSÕES</b>	<b>8</b>
<b>4. REFERÊNCIAS</b>	<b>12</b>

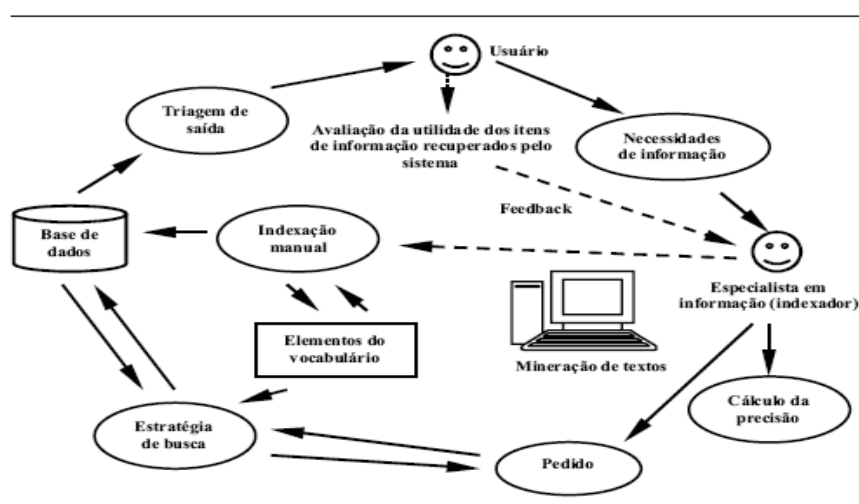
## 1. INTRODUÇÃO

A linguagem natural é a língua da maneira que os humanos comunicam e escrevem. O PLN, técnica responsável em processar esse linguajar, concilia a linguística, informática e a inteligência artificial - métodos de machine learning e deep learning. Com o uso desse recurso é possível recuperar informações escritas à natureza da língua.

Devido a interconexão entre sistemas informáticos e humanos, a quantidade de informações acessadas por computadores são cada vez maiores, assim, o processamento de linguagem natural, torna-se imprescindível para que as máquinas sejam capazes de ler e interpretá-las, a fim de entregar ao usuário uma informação enxuta e de qualidade. É possível através dessa técnica realizar sumarização, medir sentimentos e classificar as informações. Um dos exemplos mais famosos são os motores de busca que, através de pequenas frases, conseguem entregar documentos compatíveis com o que o utente deseja.

Desde os primórdios da organização social os humanos organizam as informações com intuito de facilitar a busca e a recuperação das mesmas. Um exemplo desse fato são as grandes bibliotecas, a mais antiga que se tem conhecimento foi criada em Elba, que com a criação de índices manuais que fornecem a categoria dos livros possibilitam uma melhor experiência a quem deseja encontrar algum texto.

FIGURA 1  
Fatores de influência nos resultados de busca em uma base de dados



Fonte: adaptado de Lancaster (1998).

A partir do surgimento da internet, a recuperação de informações deixou de ser vista como algo relacionado apenas às bibliotecas, mas sim algo indispensável para lidar com o volume de conhecimento humano publicado na web. Os motores de busca tem que lidar de maneira assertiva com esse extenso repositório, ao passo que tem a função de traduzir a necessidade do usuário em algo que as máquinas possam entender e entregar algo que seja condizente com as palavras chaves informadas e não documentos ruidosos. Logo, essas tecnologias devem ser programadas, testadas e melhoradas continuamente.

As informações se diferem dos dados por agregarem valor, uma vez que elas dão sentido, enquanto os dados não são capazes, por si só, de descreverem algo, nesse viés, eles são a base para a geração da informação. Portanto, apesar de não representarem a mesma coisa, ambos caminham juntos.

No presente trabalho, será utilizado um corpus(conjunto de documentos) que compreende aproximadamente 18.000 postagens de grupos de notícias classificados em 20 tópicos e divididos subconjuntos, um para treinamento e outro para teste. O conjunto “The 20 newsgroups” faz parte da biblioteca sklearn.datasets. Os arquivos serão carregados, tokenizados, filtrados, refinados e ,posteriormente, serão base para treinar um classificador que terá, em sequência, seu desempenho avaliado. O objetivo é que a classificação do documento gerada pelo sistema deve ser condizente com a real.

## 2. METODOLOGIA

O treinamento do classificador será realizado na linguagem Python por intermédio da plataforma Google Colab. Em princípio, o corpus “The 20 newsgroups” será importado do conjunto de datasets sklearn.datasets. Em seguida serão definidas as categorias de documentos que iremos utilizar, pois não iremos trabalhar com as 20 categorias existentes.

```
[276] from sklearn.datasets import fetch_20newsgroups
import numpy as np
newsgroups_train = fetch_20newsgroups(subset='train')

from pprint import pprint
pprint(list(newsgroups_train.target_names))

['alt.atheism',
 'comp.graphics',
 'comp.os.ms-windows.misc',
 'comp.sys.ibm.pc.hardware',
 'comp.sys.mac.hardware',
 'comp.windows.x',
 'misc.forsale',
 'rec.autos',
 'rec.motorcycles',
 'rec.sport.baseball',
 'rec.sport.hockey',
 'sci.crypt',
 'sci.electronics',
 'sci.med',
 'sci.space',
 'soc.religion.christian',
 'talk.politics.guns',
 'talk.politics.mideast',
 'talk.politics.misc',
 'talk.religion.misc']

[277] print(newsgroups_train.files.shape);
print(newsgroups_train.target.shape);
print(newsgroups_train.target[:10]);

(11314,)
(11314,)
[ 7  4  4  1 14 16 13  3  2  4]

[267] cats = ['sci.crypt','sci.electronics','sci.med','sci.space',]

newsgroups_train = fetch_20newsgroups(subset='train', categories=cats, random_state=0)

list(newsgroups_train.target_names)

['sci.crypt', 'sci.electronics', 'sci.med', 'sci.space']
```

fonte: autoria própria

Para prosseguir com o trabalho os textos serão tokenizados, assim, eles serão convertidos em uma representação vetorial de recursos que faz referência a frequência dos termos, além de ser filtrado para retirar as stopwords, conjunto de palavras consideradas irrelevantes para a recuperação de informações como

conjunções, e deixar o documento mais enxuto. Essas ações podem ser realizadas através tanto por meio do procedimento `TfidfVectorizer`, quanto do `CountVectorizer`. Para o próximo passo manteremos o foco no `CountVectorizer`.

```
▼ Vetorização

[270] from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer = TfidfVectorizer()
vectors = vectorizer.fit_transform(newsgroups_train.data)
vectors.shape

(2373, 38683)

[278] type(vectors)

scipy.sparse.csr.csr_matrix

[211] vectors.nnz / float(vectors.shape[0])
#em média 164 componentes diferentes de 0

164.43826380109567

▼ Tokenização

[212] from sklearn.feature_extraction.text import CountVectorizer

count_vect = CountVectorizer()
newsgroups_train_counts = count_vect.fit_transform(newsgroups_train.data)

[213] newsgroups_train_counts.shape
#tamanho da matriz

(2373, 38683)

[214] count_vect.vocabulary_.get(u'country')
#quantas vezes a palavra aparece no corpus

11807
```

fonte: autoria própria

Por meio do recurso `TfidfTransformer` em conjunto com o `CountVectorizer` a ocorrência dos termos serão transformadas em uma frequência de fato. Esse método leva em consideração não apenas quantas vezes a palavra aparece no documento, mas também a quantidade de palavras que ele contém e retorna a pontuação `tf-idf`.

Para o primeiro teste será utilizado o método multinomial Naive Bayes que considera características discretas, à exemplo, a pontuação `tf-idf`. Deste modo, serão escritos dois títulos fictícios de obras, por meio das palavras presentes o classificador irá categorizar os possíveis documentos em categorias de acordo com a `tf-idf`.

```
▼ Frequência

[309] from sklearn.feature_extraction.text import TfidfTransformer

[310] tfidf_transformer = TfidfTransformer()
      newsgroups_train_tfidf = tfidf_transformer.fit_transform(newsgroups_train_counts)
      newsgroups_train_tfidf.shape

(2373, 38683)

▼ Classificador

[311] from sklearn.naive_bayes import MultinomialNB
      classification_nb = MultinomialNB()
      classification = classification_nb.fit(newsgroups_train_tfidf, newsgroups_train.target)

[312] newsgroups_train.target_names

['sci.crypt', 'sci.electronics', 'sci.med', 'sci.space']

[313] docs_new = ['The doctor is good', 'The biggest planet']
      newsgroups_new_counts = count_vect.transform(docs_new)
      newsgroups_new_tfidf = tfidf_transformer.transform(newsgroups_new_counts)

      predicted = classification.predict(newsgroups_new_tfidf)

      for doc, category in zip(docs_new, predicted):
          print('%r => %s' % (doc, newsgroups_train.target_names[category]))

'The doctor is good' => sci.med
'The biggest planet' => sci.space
```

fonte: autoria própria

No exemplo acima os títulos 'The doctor is goo' e 'The biggest planet' foram categorizados em sci.med e sci.space, respectivamente.

Por fim será construído um pipeline, uma espécie de classificador composto que une várias fases de machine learning em um fluxo. Logo, por meio dele é possível fazer todo trabalho realizado no classificador anterior em apenas uma linha.

```
▼ Condutor

[220] from sklearn.pipeline import Pipeline
      text_clf = Pipeline([('vect', CountVectorizer()), ('tfidf', TfidfTransformer()), ('clf', MultinomialNB()),])

[221] text_clf.fit(newsgroups_train.data, newsgroups_train.target)

Pipeline(steps=[('vect', CountVectorizer()), ('tfidf', TfidfTransformer()), ('clf', LinearSVC())])
```

fonte: autoria própria

### 3. RESULTADOS E DISCUSSÕES

Os classificadores serão avaliados de diferentes modos, a fim de analisar a precisão dos modelos. O desempenho de um classificador é de suma importância, pois atender uma requisição de maneira errônea pode ocasionar diversas frustrações.

Em princípio, será analisado o desempenho do SGDC, que funciona como um classificador linear com treinamento SGD, através do Pipeline. Em sequência trocaremos apenas o método multinomial por uma máquina de vetor de suporte linear para ver se a acurácia se difere. A qualidade será medida através de quantas vezes, em média, a categoria classificada se iguala à original.

```
Desempenho do classificador

[325] import numpy as np
newsgroups_test = fetch_20newsgroups(subset='test',categories=cats, shuffle=True, random_state=42)
docs_test = newsgroups_test.data
predicted = text_clf.predict(docs_test)
np.mean(predicted == newsgroups_test.target)

0.9373020899303357

[326] from sklearn.linear_model import SGDClassifier
text_clf = Pipeline([('vect', CountVectorizer()),('tfidf', TfidfTransformer()),('clf', SGDClassifier(loss='hinge', penalty

[327] text_clf.fit(newsgroups_train.data, newsgroups_train.target)

Pipeline(steps=[('vect', CountVectorizer()), ('tfidf', TfidfTransformer()),
                ('clf',
                 SGDClassifier(alpha=0.001, max_iter=5, random_state=42,
                              tol=None))])

[328] predicted = text_clf.predict(docs_test)

[329] np.mean(predicted == newsgroups_test.target)

0.9373020899303357
```

fonte: autoria própria

É possível verificar que ambos os classificadores alcançaram 93,73% de precisão, ou seja, responderam de maneira idêntica e, para essa classificação, não existe um classificador preferível.

Nesse momento, será elaborado um relatório com as principais métricas de classificação (precisão, revocação, F-measure) para cada uma das classes utilizadas para a análise.



- Precisão: preditivo positivo, retornará um valor relacionado a quantidade de elementos relevantes recuperados, ou seja, os verdadeiros positivos em relação ao total de elementos classificados;
- Revocação: sensibilidade, refere-se a quantidade de positivos recuperados, os positivos categorizados em relação ao total de positivos reais;
- F-measure: média harmônica ponderada de precisão e revocação, mede a eficácia.

```

- Avaliação do Classificador

[228] from sklearn import metrics

- Classification report

[229] print(metrics.classification_report(newsgroups_test.target, predicted, target_names=newsgroups_test.target_names))

```

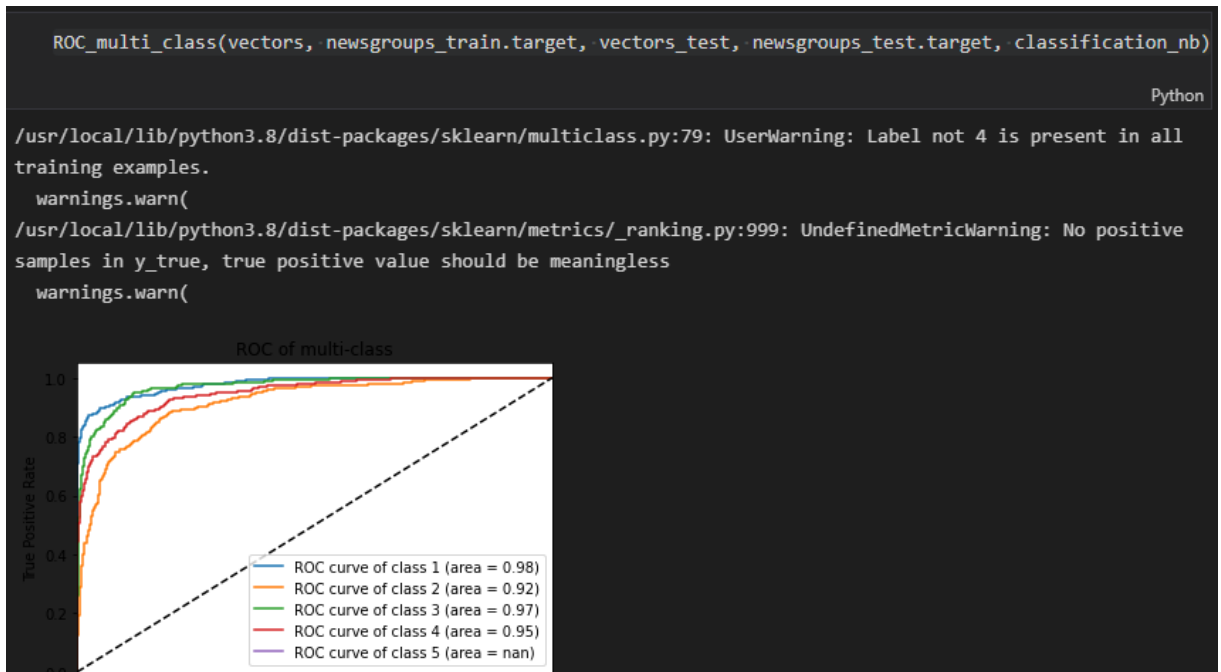
	precision	recall	f1-score	support
sci.crypt	0.98	0.94	0.96	396
sci.electronics	0.85	0.95	0.90	393
sci.med	0.96	0.90	0.93	396
sci.space	0.97	0.96	0.96	394
accuracy			0.94	1579
macro avg	0.94	0.94	0.94	1579
weighted avg	0.94	0.94	0.94	1579

fonte: autoria própria

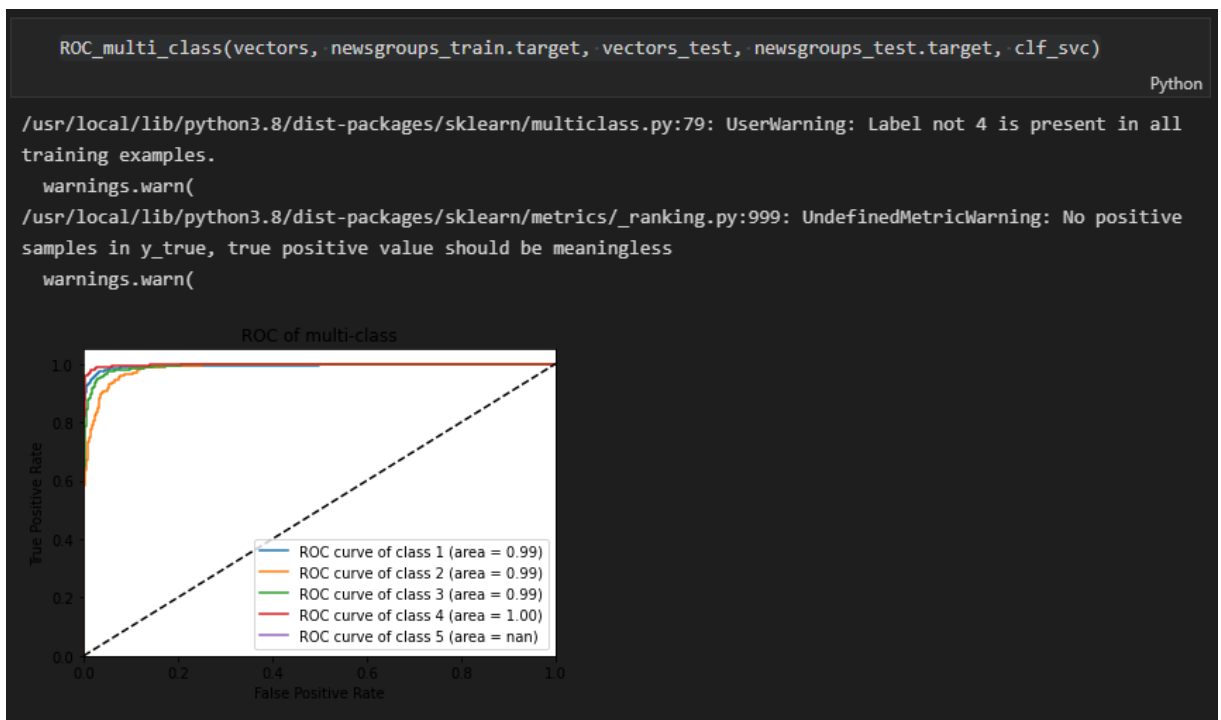
Apenas a classe sci.med apresentou uma precisão menor que 90%, o que demonstra que o classificador apresentou um alto desempenho. Além disso, a revocação foi alta em todos os casos, portanto, o número de verdadeiros recuperados foi quase idêntico ao real.

Para finalizar a avaliação da análise realizada serão plotadas curvas roque multiclasse. Os classificadores avaliados serão multinomial, linear e KNeighbors. O eixo Y dessa avaliação representa a taxa de verdadeiro positivo (TPR) resultante do classificador, enquanto o eixo X, a taxa de falso positivo(FPR). Nesse viés, o canto superior esquerdo do gráfico, com FPR zero e TPR um, é o ponto ideal.

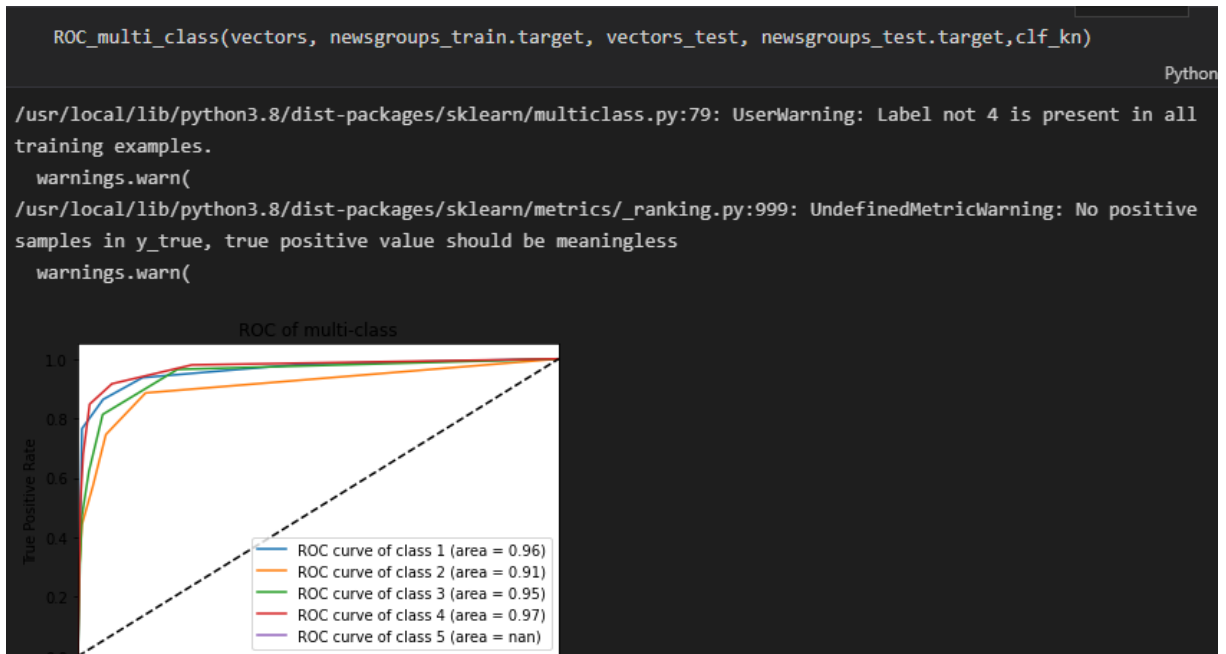
As curvas serão apresentadas todas juntas em um mesmo gráfico. Para isso será necessário utilizar a estratégia multiclasse One-vs-the-Rest, assim, em cada etapa uma determinada classe será considerada positiva e as demais negativas.



fonte: autoria própria



fonte: autoria própria



fonte: autoria própria

Os três classificadores apresentaram um bom desempenho, no entanto, a área sob a curva do modelo linear foi a maior e chegou a alcançar um valor 1. Isso pode ser explicado em razão de sua maior flexibilidade na escolha de penalidades e funções de perdas, além do fato que a classificação de vetores de suporte linear suporta um grande volume de dados e a multiplicidade de classes.

#### 4. REFERÊNCIAS

JÚNIOR, Rogério Henrique de Araújo; TARAPANOFF, KIRA. Precisão no processo de busca e recuperação da informação: uso da mineração de textos. Scielo, 2007. Disponível em: <<https://www.scielo.br/j/ci/a/wshBz99WbXhPWtD8MLQpDLd/?lang=pt#>>. Acesso em: 14 Dez. 2022.

GUIMARÃES, Leandro. Qual a diferença entre dado e informação? Entenda agora! Disponível em: <<https://www.knowsolution.com.br/diferenca-dado-e-informacao/#:~:text=Enquanto%20os%20dados%20podem%20ser, trabalho%20com%20elementos%20mais%20brutos.>>>. Acesso em: 14 Dez. 2022.

WIKIPÉDIA. Precisão e revocação. Disponível em: <[https://pt.wikipedia.org/wiki/Precis%C3%A3o\\_e\\_revoca%C3%A7%C3%A3o](https://pt.wikipedia.org/wiki/Precis%C3%A3o_e_revoca%C3%A7%C3%A3o)> Acesso em: 15 Dez. 2022.

MICROSOFT. Pipeline Classe. Disponível em: <<https://learn.microsoft.com/pt-br/python/api/azureml-pipeline-core/azureml.pipeline.core.pipeline.pipeline?view=azure-ml-py>>. Acesso em: 15 Dez. 2022.

BAEZA-YATES, R.; RIBEIRO-NETO, B. Recuperação de Informação. Bookman, 2013.. Disponível em: <<https://integrada.minhabiblioteca.com.br/#/books/9788582600498/>>. Acesso em: 16 Dec 2022.

SCIKIT-LEARN. Aprendizado de máquina em Python. Disponível em: <<https://scikit-learn.org/stable/index.html>> . Acesso em: 10 Dez. 2022.