

**Preprocessando e
extraindo dados
relevantes de
textos com tf-idf**



Igor Moura

- Sei brincar com python
- Mexo com eletrônicos
- Gosto de sofrer com computadores



@igormp

**Mas que danado é
tf-idf?**

*term frequency-inverse
document frequency* é
função matemática que
rankeia palavras dada a
sua relevância

$$TF-IDF = TF(t, d) * IDF(t, D)$$

$$TF(t, d) = \frac{\textit{count}(t, d)}{\textit{count}(*, d)}$$

$$IDF(t, D) = \log_e \left(\frac{D}{\text{amount}(t, D)} \right)$$

E como implementar mesmo?

Experimente no Colab!



Google Colaboratory

[google.com](https://colab.google.com)

```
1 !wget http://mlg.ucd.ie/files/datasets/bbc-fulltext.zip
2 !unzip bbc-fulltext.zip
3 !rm bbc-fulltext.zip*
4 %cd bbc
5 !ls
6 !rm README.TXT
```

```
1 documents = []
2 for root, dirs, files in os.walk("."):
3     for name in files:
4         if name.endswith(".txt"):
5             filepath = root + os.sep + name
6             with open(filepath, 'r', encoding='latin-1') as file:
7                 data = file.read().replace('\n', ' ').split()
8                 documents.append(data)
```

```
1 documents = [[word for word in sublist if (len(word)>1)]
2             for sublist in documents]
```

```
1 documents = [[word.lower() for word in sublist]
2             for sublist in documents]
```

```
1 symbols = "!\"#$%&()*+-./:;<=>?@[\\]^_`{|}~\'"
2 for i in symbols:
3     documents = [[''.join([char for char in word if char != i])
4                 for word in sublist]
5                 for sublist in documents]
```

```
1 nltk.download('stopwords')
2 stopwords = nltk.corpus.stopwords.words("english")
3 documents = [[word for word in sublist if word not in stopwords]
4              for sublist in documents]
```

```
1 lemmatizer = WordNetLemmatizer()
2 documents = [[lemmatizer.lemmatize(word) for word in sublist]
3              for sublist in documents]
4
5 stemmer = SnowballStemmer('english')
6 documents = [[stemmer.stem(word) for word in sublist]
7              for sublist in documents]
```

Stemming:

studies -> studi

Lemmatization:

studies -> study

- Remoção de palavras com tamanho unitário
- Deixar tudo em minúsculas
- Remover símbolos
- Remover "stopwords"
- Lemmatização seguida de stemmização

```
1 dfs = {}
2 for i in range(len(documents)):
3     for word in documents[i]:
4         try:
5             dfs[word].add(i)
6         except:
7             dfs[word] = {i}
```



```
1 for i in dfs:  
2     dfs[i] = len(dfs[i])
```

```
1 tf_idf = {}
2 for i in range(len(documents)):
3     counter = Counter(documents[i])
4     for token in np.unique(documents[i]):
5         tf = counter[token]/len(documents[i])
6         df = dfs[token]
7         idf = np.log(len(documents)/(df+1))
8         tf_idf[i, token] = tf*idf
```

```
1 (pd.DataFrame(list(tf_idf.items())))
2   .sort_values(by=[1], ascending=False)
3   .head(5))
```

	0	1
200533	(1181, rugby),)	0.716521
66307	(324, commodor)	0.492562
32093	(162, commodor)	0.492562
216404	(1297, pountney)	0.477752
91123	(451, axa)	0.458059

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 tfidf = TfidfVectorizer(sublinear_tf=True,
4                          min_df=5, norm='l2',
5                          encoding='latin-1',
6                          ngram_range=(1, 2),
7                          stop_words='english')
```

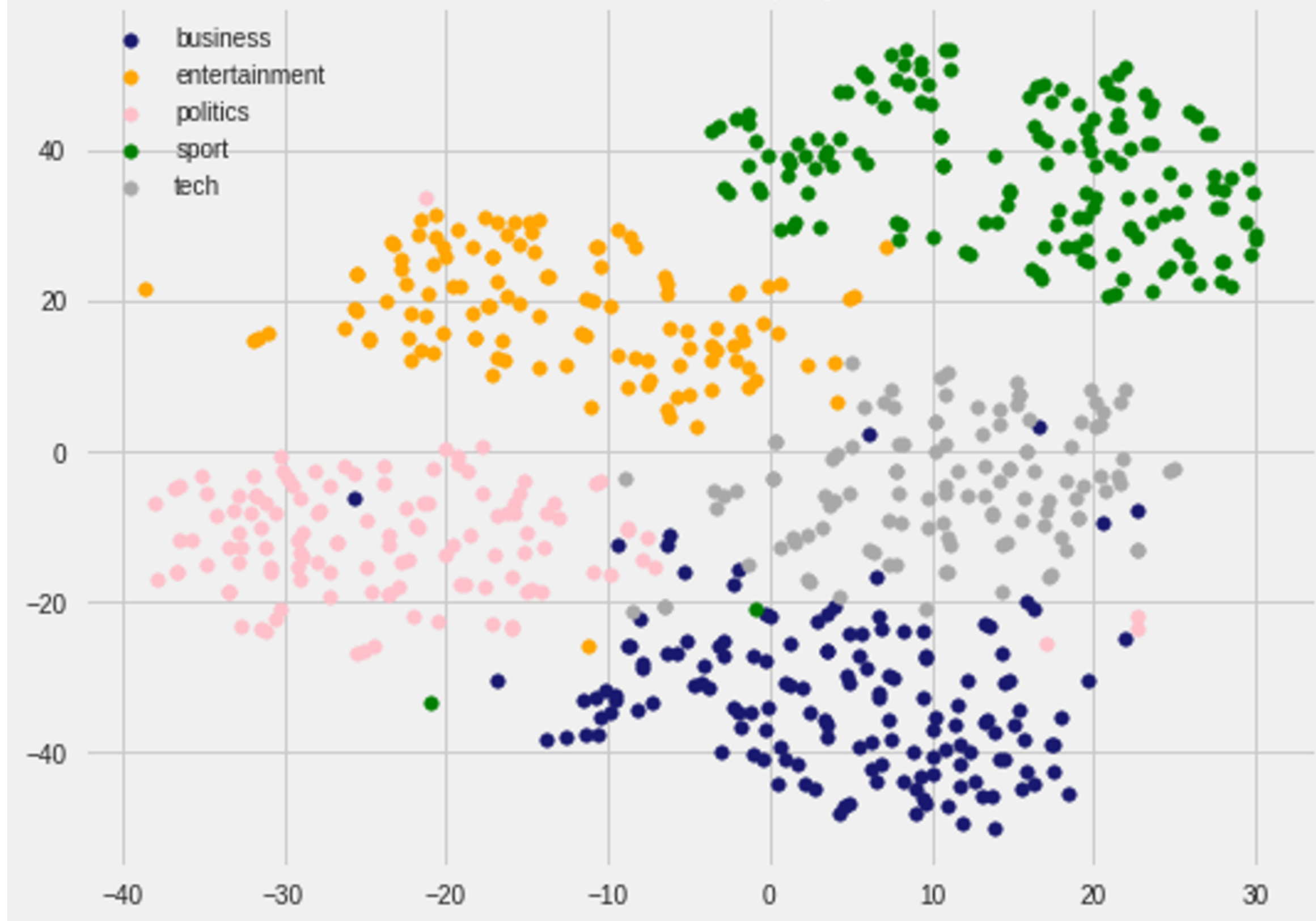
Mas qual a utilidade disso?

Ferramentas de busca

- Recuperação de informação
- Rankamento de páginas
- Extração de keywords

Natural Language Processing

tf-idf feature vector for each article, projected on 2 dimensions.



Mercado Livre



Using TensorFlow to predict product weight and dimensions

A Guest Post by Rodolfo Bonnin from the Mercado Libre Applied Machine Learning team

Medium / TensorFlow / Sep 25

TFIDF (Title)

TFIDF (Description)

DictVect (condition)

DictVect(shipping)

Dataset row

Como utilizamos na firma



MARTORELLI

ADVOGADOS

CERVEJARIA

ambev

- Petição inicial
- Ata
- Sentença
- Acórdão
- Alvarás




Obrigado!

 igormp2

 igormp

 igormp

 imp2@cin.ufpe.br

