

Análise de Sentimento em Redes Sociais Utilizando Scikit-learn+NLTK

Thais Gomes de Almeida

tga@icomp.ufam.edu.br 

DNS lab - Pesquisadores



Contextualização





Top Picks

Trending

Food

Coffee

Nightlife

Fun

Shopping



Picanha Mania

BBQ Joint \$\$\$
Centro, Manaus

Tips 181

Photos 439



Samia M.

August 4, 2016

Casa sempre cheia = atendimento um pouco confuso, mas eles são esforçados. Peça a picanha na brasa, achei melhor que a no bafo. Farofa de paçoca e queijo assado são essenciais.

Upvote

Downvote



Lumax Azevedo

February 22, 2014

Been here 5+ times

Picanha gostosa e com preço justo! O Petit Gateau e a sobremesa de morango com chantili são deliciosas :)

Upvote 7

Downvote



Luiz Branches

October 30, 2013

Ambiente agradável e amplo, ideal pra quem quer se reunir com os familiares e amigos, sem se importar em quanto vai gastar. Recomendo.

Upvote 3

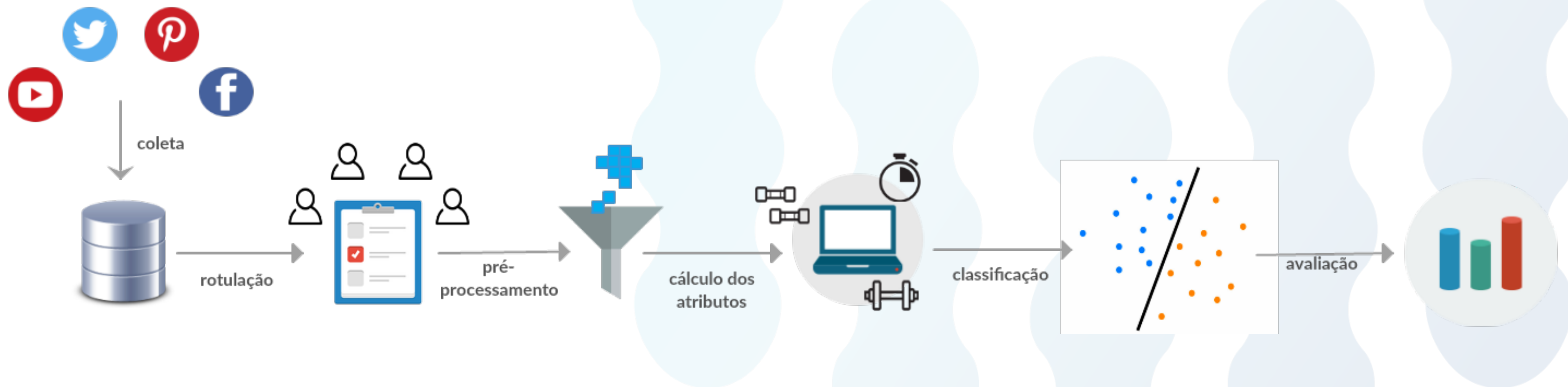
Downvote

Aprendizagem de Máquina Supervisionada

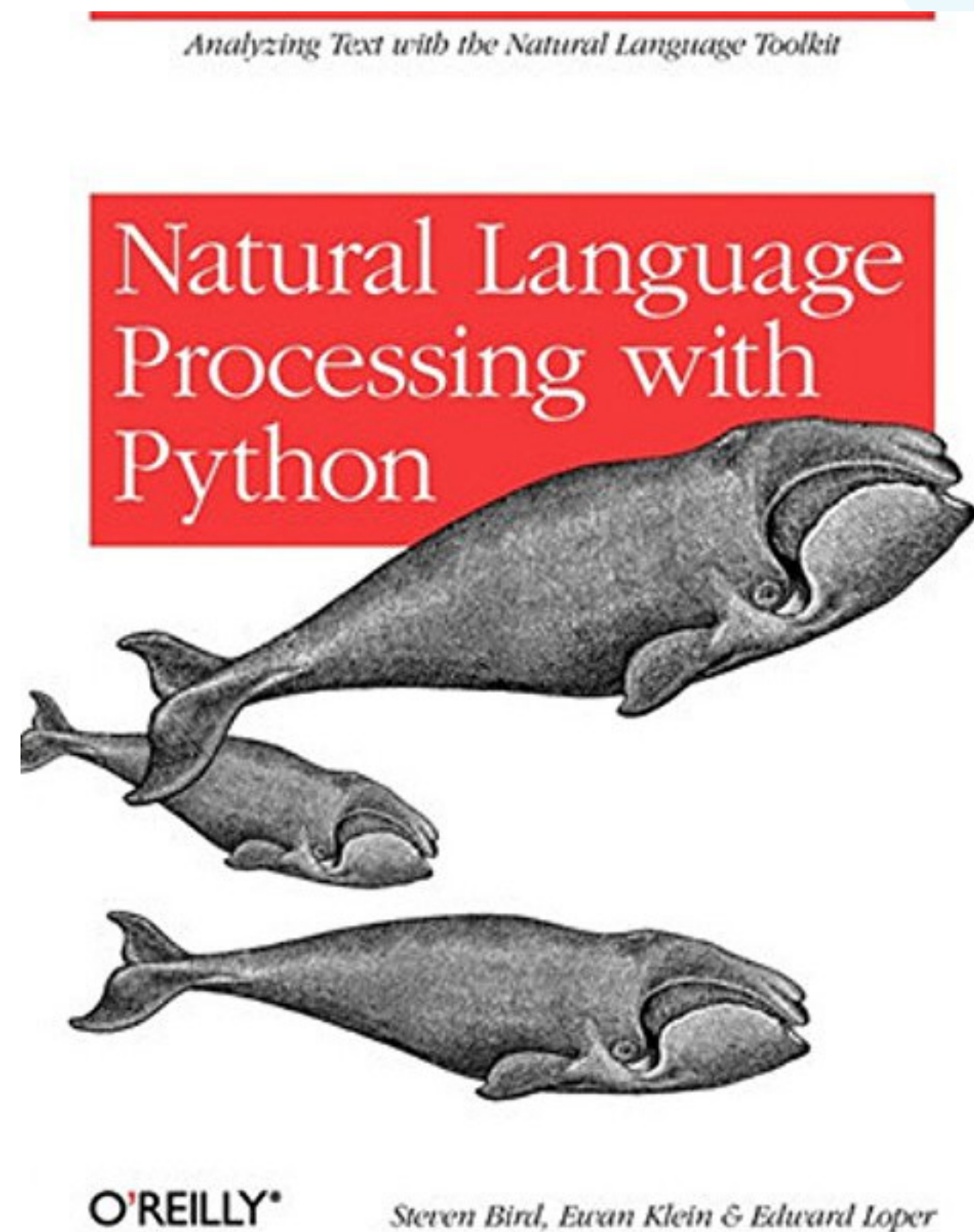


"Um conjunto de exemplos de **treinamento** com as **respostas corretas** (alvos) é fornecido e, com base neste conjunto de treinamento, o algoritmo geral consegue responder corretamente a todas as entradas possíveis." (Herbrich e Graepel, 2014)

Arquitetura Clássica de Aprendizagem Supervisionada



Natural Language Toolkit - NLTK



- Atualmente na **versão 3.4.2**;
- Torna ágil operações como *tokenization*, *stemming*, *lemmatization*;
- Possui implementações de métodos supervisionados e de Reconhecimento de Entidade Nomeada;

NLTK: Pré-Processamento

```
1 list_docs = []
2 list_label = []
3
4 with open('dataset_foursquare_tips.csv') as csvfile:
5     dataset = csv.reader(csvfile, delimiter=',')
6     for line in dataset:
7         # pre-process considering words
8         tips_text = remove_ulr_corpus(line[0])
9         tips_text = remove_mention_corpus(tips_text)
10        tips_text = tokenize_words_corpus(tips_text)
11        tips_text = remove_stopwords_corpus(tips_text)
12        document = ' '.join(tips_text)
13
14        list_docs.append(document)
15        list_label.append(line[1])
```

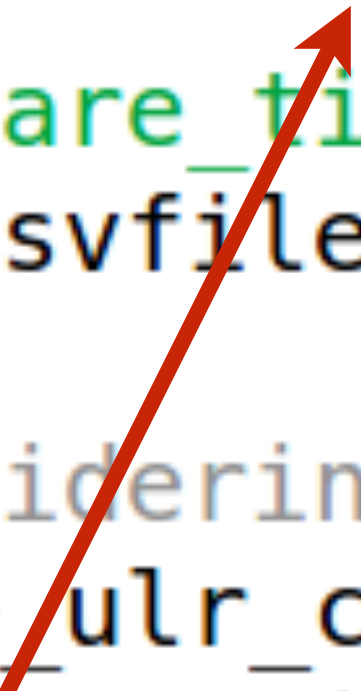

NLTK: Pré-Processamento

```
1 list_docs = []
2 list_label = re.sub(r'^https?:\/\/\/*.*[\r\n]*', '', text)
3
4 with open('dataset_foursquare_tips.csv') as csvfile:
5     dataset = csv.reader(csvfile, delimiter=',')
6     for line in dataset:
7         # pre-process considering words
8         tips_text = remove_ulr_corpus(line[0])
9         tips_text = remove_mention_corpus(tips_text)
10        tips_text = tokenize_words_corpus(tips_text)
11        tips_text = remove_stopwords_corpus(tips_text)
12        document = ' '.join(tips_text)
13
14        list_docs.append(document)
15        list_label.append(line[1])
```

NLTK: Pré-Processamento

```
1 list_docs = []
2 list_label = []
3
4 with open('dataset_foursquare_tips.csv') as csvfile:
5     dataset = csv.reader(csvfile, delimiter=',')
6     for line in dataset:
7         # pre-process considering words
8         tips_text = remove_url_corpus(line[0])
9         tips_text = remove_mention_corpus(tips_text)
10        tips_text = tokenize_words_corpus(tips_text)
11        tips_text = remove_stopwords_corpus(tips_text)
12        document = ' '.join(tips_text)
13
14        list_docs.append(document)
15        list_label.append(line[1])
```

`re.sub(r'@\w+ ?', '', text)`



NLTK: Pré-Processamento

```
1 from nltk.tokenize import RegexpTokenizer
2 from nltk.corpus import stopwords
3
4 def tokenize_words_corpus(text):
5     tokenizer = RegexpTokenizer(r'\w+')
6     tokens = tokenizer.tokenize(text)
7     return tokens
8
9
10 from nltk.stem import RSLPStemmer
11 from nltk.tokenize import word_tokenize
12
13 def tokenize_words_corpus(text):
14     tokens = word_tokenize(text)
15     porter = RSLPStemmer()
16     tokens = [porter.stem(t) for t in tokens]
17     tokens = [t for t in tokens
18               if len(t) > 2 and not t.isdigit()]
19     return tokens
```

NLTK: Pré-Processamento

```
1 from nltk.corpus import stopwords
2
3 def remove_stopwords_corpus(tip_text):
4     content = []
5     for word in tip_text:
6         if word.lower().strip() not in stopwords.words('portuguese'):
7             content.append(word.lower().strip())
8     return content
9
```


NLTK: Pré-Processamento

```
1 list_docs = []
2 list_label = []
3
4 with open('dataset_foursquare_tips.csv') as csvfile:
5     dataset = csv.reader(csvfile, delimiter=',')
6     for line in dataset:
7         # pre-process considering words
8         tips_text = remove_ulr_corpus(line[0])
9         tips_text = remove_mention_corpus(tips_text)
10        tips_text = tokenize_words_corpus(tips_text)
11        tips_text = remove_stopwords_corpus(tips_text)
12        document = ' '.join(tips_text)
13
14        list_docs.append(document)
15        list_label.append(line[1])
```

Scikit-learn



- Atualmente na **versão 0.19** (julho 17);
- Está integrado com o **SciPy + NumPy**;
- Contém métodos de classificação, regressão, agrupamento...e muito mais!;

Scikit-learn

```
1 from sklearn.model_selection import StratifiedKFold, cross_val_predict
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics import recall_score, precision_score, f1_score
5
6 #data representation
7 vectorizer = TfidfVectorizer(use_idf=True, sublinear_tf=False)
8 X = vectorizer.fit_transform(list_docs)
9 Y = list_label
10
11 #text classifier
12 knn = KNeighborsClassifier(n_neighbors=5)
13
14 #training
15 stratified_fold = StratifiedKFold(n_splits=10, shuffle=True)
16 predict_list = cross_val_predict(knn, X, Y, cv = stratified_fold)
17
18 #evaluate training performance
19 precision = precision_score(Y, predict_list, average=None)
20 recall = recall_score(Y, predict_list, average=None)
21 f1 = f1_score(Y, predict_list, average=None)
```

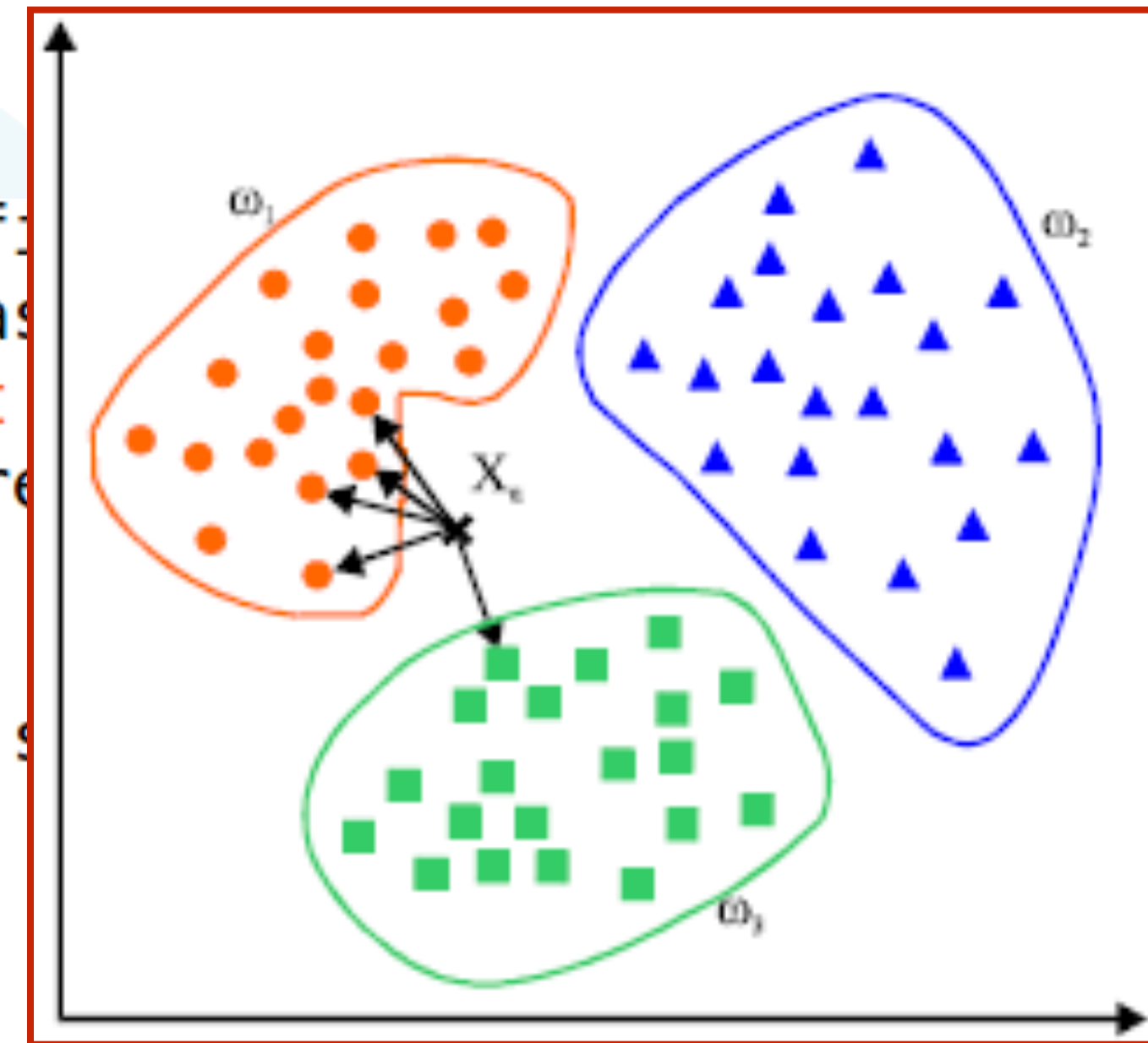
Scikit-learn

```
1 from sklearn.model_selection
2 from sklearn.neighbors import
3 from sklearn.feature_extract
4 from sklearn.metrics import
5
6 #data representation
7 vectorizer = TfidfVectorizer(use_idf=True, sublinear_tf=False)
8 X = vectorizer.fit_transform(list_docs)
9 Y = list_label
10
11 #text classifier
12 knn = KNeighborsClassifier(n_neighbors=5)
13
14 #training
15 stratified_fold = StratifiedKFold(n_splits=10, shuffle=True)
16 predict_list = cross_val_predict(knn, X, Y, cv = stratified_fold)
17
18 #evaluate training performance
19 precision = precision_score(Y, predict_list, average=None)
20 recall = recall_score(Y, predict_list, average=None)
21 f1 = f1_score(Y, predict_list, average=None)
```

$$w_{i,j} = \begin{cases} (1 + \log f_{i,j}) \times \log \frac{N}{n_i} & \text{se } f_{i,j} > 0 \\ 0 & \text{caso contrário} \end{cases}$$

Scikit-learn

```
1 from sklearn.model_selection import StratifiedKFold
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics import recall_score, precision_score
5
6 #data representation
7 vectorizer = TfidfVectorizer(use_idf=True, smooth_idf=0.5)
8 X = vectorizer.fit_transform(list_docs)
9 Y = list_label
10
11 #text classifier
12 knn = KNeighborsClassifier(n_neighbors=5)
13
14 #training
15 stratified_fold = StratifiedKFold(n_splits=10, shuffle=True)
16 predict_list = cross_val_predict(knn, X, Y, cv = stratified_fold)
17
18 #evaluate training performance
19 precision = precision_score(Y, predict_list, average=None)
20 recall = recall_score(Y, predict_list, average=None)
21 f1 = f1_score(Y, predict_list, average=None)
```



Scikit-learn

```
1 from sklearn.model_selection import StratifiedKFold, cross_val_predict
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics import recall_score, precision_score, f1_score
5
6 #data representation
7 vectorizer = TfidfVectorizer(use_idf=True, sublinear_tf=False)
8 X = vectorizer.fit_transform(list_docs)
9 Y = list_label
10
11 #text classifier
12 knn = KNeighborsClassifier(n_neighbors=5)
13
14 #training
15 stratified_fold = StratifiedKFold(n_splits=10, shuffle=True)
16 predict_list = cross_val_predict(knn, X, Y, cv = stratified_fold)
17
18 #evaluate training performance
19 precision = precision_score(Y, predict_list, average=None)
20 recall = recall_score(Y, predict_list, average=None)
21 f1 = f1_score(Y, predict_list, average=None)
```


Scikit-learn

```
1 from sklearn.model_selection import StratifiedKFold, cross_val_predict
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.metrics import recall_score, precision_score, f1_score
5
6 #data representation
7 vectorizer = TfidfVectorizer(use_idf=True, sublinear_tf=False)
8 X = vectorizer.fit_transform(list_docs)
9 Y = list_label
10
11 #text classifier
12 knn = KNeighborsClassifier(n_neighbors=5)
13
14 #training
15 stratified_fold = StratifiedKFold(n_splits=10, shuffle=True)
16 predict_list = cross_val_predict(knn, X, Y, cv = stratified_fold)
17
18 #evaluate training performance
19 precision = precision_score(Y, predict_list, average=None)
20 recall = recall_score(Y, predict_list, average=None)
21 f1 = f1_score(Y, predict_list, average=None)
```

Considerações Finais



NLTK



Referências Bibliográficas

- Baeza-Yates, Ricardo, and Berthier Ribeiro-Neto. Recuperação de Informação-: Conceitos e Tecnologia das Máquinas de Busca. Bookman Editora, 2013.
- Marsland, Stephen. Machine learning: an algorithmic perspective. CRC press, 2015.
- Bird, Steven, Edward Loper and Ewan Klein (2009), Natural Language Processing with Python. O'Reilly Media Inc;



Obrigada!



Thais Gomes de Almeida
tga@icomp.ufam.edu.br ✉