

Processamento de Linguagem Natural: Vetorização de Texto com Python

Processamento de Linguagem Natural: Vetorização de Texto com Python

Giseldo Neo

Alana Neo

Preparação e revisão Giseldo Neo
Diagramação Alana Neo

©2024
versão 0.1

E-mail: giseldo@gmail.com

Todos os direitos reservados.
Nenhuma parte desta publicação
poderá ser armazenada ou reproduzida
por qualquer meio sem a autorização
por escrito dos autores.

Para minhas filhas.

Informações Adicionais

Este Livro está em uma versão Beta e em constante evolução.

O código fonte dos exemplos no livro pode ser encontrada em:

<https://giseldo.github.io>

Se você encontrou algum erro, deseja enviar alguma sugestão ou está com alguma dúvida, envie um e-mail para giseldo@gmail.com

Prefácio

O processamento de linguagem natural (PLN) tem se tornado cada vez mais relevante na era da informação. Este livro oferece uma introdução concisa, porém abrangente, ao processo de vetorização de texto, uma técnica fundamental em PLN. Espero que este livro seja um ponto de partida útil para suas futuras explorações no mundo do processamento de linguagem natural.

Boa leitura!

Giseldo Neo

Sumário

1	Introdução à Vetorização de Texto	13
1.1	Conceito de Vetorização de Texto	13
1.2	Importância na Análise de Dados e Aprendizado de Máquina	14
1.3	Aplicações Práticas	16
2	Preparação do Ambiente	21
2.1	Instalação do Python	21
2.1.1	Instalando o Python	21
2.2	Instalação de Bibliotecas Necessárias	22
2.2.1	Instalando Bibliotecas com pip	22
2.2.2	Exemplo em Python: Verificando Instalações . . .	22
2.3	Introdução ao Jupyter Notebook	23
2.3.1	Instalando o Jupyter Notebook	23
2.3.2	Exemplo em Python: Primeiros Passos no Jupyter .	24
3	Pré-processamento de Texto	27
3.1	Limpeza de Texto	27
3.1.1	Exemplo em Python: Limpeza de Texto	27
3.2	Tokenização	29
3.2.1	Exemplo em Python: Tokenização	29
3.3	Lematização e Stemming	30
3.3.1	Exemplo em Python: Lematização e Stemming . .	30

4	Vetorização de Texto	35
4.1	Bag of Words (BoW)	35
4.1.1	Exemplo em Python: Bag of Words	37
4.2	TF-IDF	38
4.2.1	Exemplo em Python: TF-IDF	39
4.3	Word Embeddings	40
4.3.1	Exemplo em Python: Word2Vec	41
5	Modelos Avançados de Vetorização	45
5.1	Embeddings Contextuais	45
5.1.1	Exemplo em Python: Usando BERT para Vetorização	46
5.2	Análise de Sentimento com Embeddings	47
5.2.1	Exemplo em Python: Classificação de Sentimento com BERT	48
5.3	Redução de Dimensionalidade	49
5.3.1	Exemplo em Python: Redução de Dimensionalidade com PCA	51
5.3.2	Exemplo em Python: Redução de Dimensionalidade com t-SNE	53
6	Aplicações Práticas	57
6.1	Classificação de Texto	57
6.1.1	Exemplo em Python: Classificação de Texto com TF-IDF e Naive Bayes	60
6.2	Agrupamento de Documentos	61
6.2.1	Exemplo em Python: Agrupamento com k-means	62
6.3	Detecção de Tópicos	63
6.3.1	Exemplo em Python: Detecção de Tópicos com LDA	64
7	Estudo de Caso	69
7.1	Análise de Reviews de Produtos	69
7.1.1	Exemplo em Python: Análise de Sentimentos em Reviews de Produtos	71
7.2	Processamento de Tweets	72
7.2.1	Exemplo em Python: Análise de Sentimentos em Tweets	73

SUMÁRIO

7.3	Análise de Notícias	74
7.3.1	Exemplo em Python: Modelagem de Tópicos em Notícias	75
8	Próximos Passos	79
8.1	Futuras Direções no Campo de Vetorização de Texto	79
8.1.1	Modelos de Linguagem de Grande Escala	79
8.1.2	Multimodalidade	79
8.1.3	Vetorização de Texto em Tempo Real	80
8.2	Leitura Recomendada	80
8.3	Próximos Passos	80

Introdução à Vetorização de Texto

1.1 Conceito de Vetorização de Texto

O processamento de linguagem natural (PLN) é um campo interdisciplinar que se concentra em permitir que computadores compreendam, interpretem e gerem linguagem humana. Uma etapa fundamental no PLN é a vetorização de texto, que consiste em transformar texto, um dado textual, em uma representação numérica que possa ser processada por algoritmos de machine learning.

Algoritmos de machine learning, como redes neurais e algoritmos de clustering, operam em espaços vetoriais. Ao representar palavras e documentos como vetores numéricos, podemos aplicar essas ferramentas para realizar tarefas como classificação, clustering e geração de texto.

Podemos pensar na vetorização de texto como uma forma de mapear palavras e documentos para um espaço vetorial multidimensional. Cada dimensão desse espaço representa uma característica particular do texto, como a frequência de uma palavra, o contexto em que ela aparece ou sua relação semântica com outras palavras.

Por exemplo, a palavra “cachorro” poderia ser representada por um vetor numérico de 70 dimensões, onde cada dimensão corresponde a uma característica específica, como a presença da palavra em um contexto relacionado a animais, a sua similaridade com a palavra "cão", etc. Ao transformar texto em vetores, podemos aplicar algoritmos de machine learning

para extrair informações valiosas e realizar tarefas complexas.

A vetorização de texto é o processo de transformar texto em uma representação numérica que pode ser utilizada por algoritmos de aprendizado de máquina.

1.2 Importância na Análise de Dados e Aprendizado de Máquina

Na era dos dados digitais, o texto é uma das formas mais abundantes de dados não estruturados. Analisar e extrair informações úteis de texto requer sua conversão para uma forma numérica, permitindo a aplicação de métodos estatísticos e de aprendizado de máquina. Uma das principais razões para a vetorização é que a maioria dos modelos matemáticos e algoritmos de machine learning requerem números como entrada, em vez de texto bruto. Ela é essencial para tarefas como classificação de texto, análise de sentimentos, detecção de tópicos e muitos outros.

Para ilustrar a importância da vetorização, vamos começar com um exemplo simples de contagem de palavras em um conjunto de documentos. Este exemplo utiliza a biblioteca `CountVectorizer` do `Scikit-learn`.

Dependências necessárias

```
! pip install scikit-learn
```

Código Python

```
1 from sklearn.feature_extraction.text import
    CountVectorizer
2
3 # Exemplo de documentos
4 documentos = [
5     "Texto de exemplo para vetorização",
6     "Outro exemplo de texto",
7     "Vetorização de texto é essencial"
8 ]
9
10 # Criação do vetor de contagem
11 vectorizer = CountVectorizer()
12 X = vectorizer.fit_transform(documentos)
13
14 # Exibição do vetor de características
15 print("Vetor de características:\n", vectorizer.
        get_feature_names_out())
16 print("Matriz BoW:\n", X.toarray())
```

Saída do Console

```
1 Vetor de características:
2 ['de' 'essencial' 'exemplo' 'outro' 'para' 'texto' '
    vetorização']
3 Matriz BoW:
4 [[1 0 1 0 1 1 1]
5  [1 0 1 1 0 1 0]
6  [1 1 0 0 0 1 1]]
```

Neste exemplo, as frases são transformadas em uma matriz de contagem de palavras (*Bag of Words*), onde cada linha representa um documento e cada coluna representa uma palavra do vocabulário.

1.3 Aplicações Práticas

A vetorização de texto encontra aplicação em diversas áreas, abrangendo múltiplos domínios do conhecimento. No campo da mineração de dados, por exemplo, a vetorização de texto é utilizada para transformar grandes volumes de dados textuais em formas estruturadas que permitem a análise quantitativa. Nos motores de busca, essa técnica possibilita que as palavras-chaves nas pesquisas dos usuários sejam comparadas eficientemente com o conteúdo indexado, proporcionando resultados mais relevantes.

Na aprendizagem de máquina, a vetorização de texto é essencial para alimentar algoritmos com entradas numéricas derivadas de textos, permitindo o treinamento de modelos para tarefas como classificação de sentimentos, detecção de spam e tradução automática. Além disso, em áreas como a análise de redes sociais, a vetorização de texto facilita a identificação de tendências e padrões ao converter postagens e comentários em matrizes numéricas, possibilitando a análise estatística. A versatilidade da vetorização de texto abre caminhos para sua aplicação em variados contextos, auxiliando na extração de insights valiosos a partir de dados textuais.

Algumas aplicações práticas

- Classificação de texto: Determinar a categoria de um documento (ex.: spam ou não-spam).
- Análise de sentimentos: Identificar a polaridade de opiniões em textos (ex.: positiva, negativa ou neutra). Por exemplo, a detecção de tendências e sentimentos em plataformas como Twitter e Facebook [16]
- Sistemas de recomendação: Sugerir produtos ou conteúdos com base em descrições textuais.
- Busca e recuperação de informação: Melhorar a relevância dos resultados de busca analisando o conteúdo dos documentos. Motores de busca: Utilizam vetorização de texto para indexar páginas web e retornar resultados relevantes para as consultas dos usuários [13]
- Assistentes virtuais: Ferramentas como Alexa, Siri, ChatGPT e Cloude utilizam técnicas de vetorização para entender e responder a comandos de voz [11].

Em resumo, a vetorização de texto é um passo essencial para qualquer tarefa de processamento de linguagem natural (PLN). Compreender as diferentes técnicas de vetorização e sua aplicação prática permite a construção de modelos mais precisos e eficientes para análise de dados textuais.

Exercícios

Versão on-line destes exercícios

<https://forms.gle/fz6gmAz5GYE2jHJF6>

1. O que é vetorização de texto?

- (a) Processo de transformar números em texto.
- (b) Processo de transformar texto em uma representação numérica.
- (c) Técnica de compressão de arquivos de texto.
- (d) Método para gerar texto automaticamente.

2. Qual das seguintes opções é uma aplicação da vetorização de texto?

- (a) Envio de textos pela internet.
- (b) Compressão de arquivos de texto.
- (c) Classificação de texto.
- (d) Edição de documentos de texto.

3. Por que a vetorização de texto é importante no aprendizado de máquina?

- (a) Porque algoritmos de aprendizado de máquina funcionam diretamente com texto bruto.
- (b) Porque algoritmos de aprendizado de máquina requerem dados numéricos para processamento.
- (c) Porque a vetorização é a única maneira de analisar grandes volumes de texto.
- (d) Porque a vetorização simplifica a criação de modelos visuais.

4. Qual é uma das principais razões para utilizar a vetorização de texto em projetos de aprendizado de máquina?

- (a) Para melhorar a legibilidade dos textos.

- (b) Para permitir que algoritmos de aprendizado de máquina processem dados textuais, que precisam ser convertidos em números.
- (c) Para aumentar o tamanho do corpus de texto.
- (d) Para corrigir erros ortográficos automaticamente.

5. Qual das seguintes etapas não faz parte do pré-processamento de texto?

- (a) Limpeza de texto.
- (b) Tokenização.
- (c) Treinamento de modelo.
- (d) Remoção de stopwords.

Capítulo 2

Preparação do Ambiente

Neste capítulo, abordaremos como preparar o ambiente necessário para trabalhar com vetorização de texto em Python. Isso inclui a instalação do Python e das bibliotecas necessárias, além de uma breve introdução ao uso do Jupyter Notebook.

2.1 Instalação do Python

Para começar a trabalhar com vetorização de texto, é essencial ter o Python instalado. Python é uma linguagem de programação amplamente utilizada para análise de dados e aprendizado de máquina devido à sua simplicidade e a vasta gama de bibliotecas disponíveis.

2.1.1 Instalando o Python

Para instalar o Python, siga as instruções abaixo:

- No Windows, baixe o instalador do site oficial do Python (<https://www.python.org/>) e siga as instruções do instalador.
- No macOS, você pode usar o Homebrew para instalar o Python executando o comando `brew install python`.

- No Linux, o Python geralmente já está instalado, mas você pode atualizá-lo usando o gerenciador de pacotes da sua distribuição.

2.2 Instalação de Bibliotecas Necessárias

Uma vez que o Python esteja instalado, precisamos instalar algumas bibliotecas que são fundamentais para a vetorização de texto. Entre as principais estão NumPy, Pandas, Scikit-learn, NLTK e SpaCy.

2.2.1 Instalando Bibliotecas com `pip`

O `pip` é o gerenciador de pacotes padrão do Python. Você pode instalar as bibliotecas necessárias usando o seguinte comando:

Dependências
<pre>pip install numpy pandas scikit-learn nltk spacy</pre>

2.2.2 Exemplo em Python: Verificando Instalações

Após instalar as bibliotecas, é importante verificar se elas foram instaladas corretamente:

Código Python

```
1 import numpy as np
2 import pandas as pd
3 import sklearn
4 import nltk
5 import spacy
6
7 print("NumPy version:", np.__version__)
8 print("Pandas version:", pd.__version__)
9 print("Scikit-learn version:", sklearn.__version__)
10 print("NLTK version:", nltk.__version__)
11 print("SpaCy version:", spacy.__version__)
```

Saída do Console

```
1 NumPy version: 1.26.4
2 Pandas version: 2.2.2
3 Scikit-learn version: 1.5.1
4 NLTK version: 3.9.1
5 SpaCy version: 3.7.6
```

Este código importará as bibliotecas e exibirá suas versões, garantindo que todas estejam corretamente instaladas.

2.3 Introdução ao Jupyter Notebook

O Jupyter Notebook é uma ferramenta poderosa para o desenvolvimento de scripts em Python, permitindo a combinação de código, texto, visualizações e resultados em um único documento.

2.3.1 Instalando o Jupyter Notebook

Você pode instalar o Jupyter Notebook usando o `pip`:

Código Python

```
1 pip install jupyterlab
```

Para iniciar o Jupyter Notebook, execute o seguinte comando no terminal:

No terminal

```
1 jupyter notebook
```

Isso abrirá o Jupyter Notebook no seu navegador padrão, permitindo que você comece a escrever e executar código Python de maneira interativa.

2.3.2 Exemplo em Python: Primeiros Passos no Jupyter

Um exemplo simples de uso do Jupyter Notebook seria a criação de uma célula de código para calcular a soma de dois números:

Código Python

```
1 a = 10
2 b = 20
3 print("A soma de a e b é:", a + b)
```

Este exemplo demonstra a simplicidade e interatividade que o Jupyter Notebook oferece, permitindo que você execute código Python célula por célula e veja os resultados imediatamente.

Em resumo, configuramos o ambiente necessário para trabalhar com vetorização de texto em Python. Com Python e as principais bibliotecas instaladas, além da configuração do Jupyter Notebook, você está pronto para explorar e aplicar técnicas de vetorização de texto.

Exercícios

Versão on-line destes exercícios

<https://forms.gle/wXiSykwHlQRdLovx5>

1. Qual é o objetivo principal da preparação do ambiente para vetorização de texto em Python?
 - (a) Instalar editores de texto avançados.
 - (b) Configurar o ambiente de desenvolvimento com Python e bibliotecas necessárias.
 - (c) Criar um ambiente gráfico para visualização de dados.
 - (d) Desenvolver interfaces de usuário.
2. Qual é o objetivo da biblioteca Scikit-learn?
 - (a) Visualizar gráficos e imagens.
 - (b) Manipular dados em planilhas.
 - (c) Fornecer ferramentas para aprendizado de máquina, incluindo vetorização de texto.
 - (d) Realizar cálculos matemáticos avançados.
3. Para que serve a biblioteca NLTK em um projeto de vetorização de texto?
 - (a) Para visualização de gráficos.
 - (b) Para manipulação de grandes volumes de dados numéricos.
 - (c) Para processamento e análise de texto natural.
 - (d) Para criação de modelos de aprendizado profundo.
4. Qual ferramenta permite a criação de documentos interativos combinando código, texto e visualizações?
 - (a) Ganymed.
 - (b) Jupyter Notebook.

- (c) Ms Paint.
 - (d) SQL Server.
5. Qual comando é utilizado para instalar a biblioteca Scikit-learn usando pip?
- (a) `pip install numpy`
 - (b) `pip install pandas`
 - (c) `pip install scikit-learn`
 - (d) `pip install matplotlib`

Pré-processamento de Texto

O pré-processamento de texto é uma etapa crucial na análise de dados textuais e no aprendizado de máquina. Antes de aplicar técnicas de vetorização, é essencial transformar e limpar os dados de texto para que possam ser processados eficientemente pelos algoritmos. Neste capítulo, abordaremos as principais técnicas de pré-processamento de texto, incluindo limpeza, tokenização, lematização e stemming.

3.1 Limpeza de Texto

A limpeza de texto envolve a remoção de elementos indesejados, como stopwords, pontuação, números e caracteres especiais, que não contribuem para a análise. Abaixo está um exemplo em Python de como realizar a limpeza básica de texto usando a biblioteca *re* para expressões regulares e *NLTK* para remoção de stopwords.

3.1.1 Exemplo em Python: Limpeza de Texto

O código a seguir remove as stopwords e pontuação de determinado texto.

Código Python

```
1 import re
2 import nltk
3 from nltk.corpus import stopwords
4
5 # Certifique-se de baixar as stopwords
6 nltk.download('stopwords')
7
8 # Exemplo de texto
9 texto = 'O processo de vetorização \nde texto envolve v
    árias etapas, como a limpeza do texto!'
10
11 # Convertendo para minúsculas
12 texto = texto.lower()
13
14 # Removendo pontuação e caracteres especiais
15 texto = re.sub(r'[^\\w\\s]', '', texto)
16 print('Texto sem pontuação:', texto)
17
18 # Removendo stopwords
19 stop_words = set(stopwords.words('portuguese'))
20 texto_limpo = ' '.join([palavra for palavra in texto.
    split() if palavra not in stop_words])
21 print('Texto limpo:', texto_limpo)
```

Saída do Console

```
1 Texto sem pontuação: o processo de vetorização de
    texto envolve várias etapas como a limpeza do texto
2 Texto limpo: processo vetorização texto envolve várias
    etapas limpeza texto
```

Este código converte o texto para minúsculas, remove pontuação e palavras que não trazem sentido semântico, tais como stopwords, resultando em uma versão limpa do texto pronta para vetorização.

3.2 Tokenização

A tokenização é o processo de dividir o texto em unidades menores, como palavras ou frases. Essas unidades são chamadas de tokens. A tokenização pode ser feita de várias formas, dependendo da granularidade desejada. A seguir, mostramos como tokenizar um texto usando a biblioteca *NLTK*.

3.2.1 Exemplo em Python: Tokenização

Código Python

```
1 import nltk
2
3 # Certifique-se de baixar o tokenizer
4 nltk.download('punkt')
5
6 from nltk.tokenize import word_tokenize
7
8 # Exemplo de texto
9 texto = 'Tokenização é o processo de dividir o texto
10         em palavras ou frases.'
11
12 # Tokenizando o texto
13 tokens = word_tokenize(texto)
14
15 print('Tokens:', tokens)
```

Console

```
1 Tokens: ['Tokenização', 'é', 'o', 'processo', 'de', 'dividir', 'o', 'texto', 'em', 'palavras', 'ou', 'frases', '.']
```

Neste exemplo, o texto é dividido em tokens individuais, que podem ser utilizados para análise posterior, como vetorização.

3.3 Lematização e Stemming

Lematização e stemming são técnicas para reduzir as palavras às suas formas base ou raízes. A lematização considera o contexto e reduz as palavras ao seu lema, enquanto o stemming simplesmente corta os sufixos para encontrar a raiz.

3.3.1 Exemplo em Python: Lematização e Stemming

Este código demonstra como aplicar stemming e lematização em um conjunto de palavras, resultando em suas formas raiz e lema, respectivamente.

Código Python

```
1 from nltk.stem import PorterStemmer, WordNetLemmatizer
2
3 # Certifique-se de baixar o WordNet
4 nltk.download('wordnet')
5 nltk.download('omw-1.4')
6
7 # Exemplo de texto
8 palavras = ['running', 'jumps', 'easily', 'fairly']
9
10 # Inicializando Stemmer e Lemmatizer
11 stemmer = PorterStemmer()
12 lemmatizer = WordNetLemmatizer()
13
14 # Aplicando Stemming e Lematização
15 stems = [stemmer.stem(palavra) for palavra in palavras]
16 lemmas = [lemmatizer.lemmatize(palavra) for palavra in
17           palavras]
18 print('Stemming:', stems)
19 print('Lematização:', lemmas)
```

Saída do Console

```
1 Stemming: ['run', 'jump', 'easili', 'fairli']  
2 Lematização: ['running', 'jump', 'easily', 'fairly']
```

Em resumo, o pré-processamento de texto é uma etapa essencial para garantir que os dados textuais estejam em uma forma adequada para a vetorização e para o aprendizado de máquina. Técnicas como limpeza, tokenização, lematização e stemming são fundamentais para preparar o texto para análise.

Exercícios

Versão on-line destes exercícios

<https://forms.gle/FFGEvmilzmrQ5Npk9>

1. Qual é o objetivo principal do pré-processamento de texto?

- (a) Melhorar a legibilidade do texto.
- (b) Transformar o texto em uma forma que possa ser processada por algoritmos de aprendizado de máquina.
- (c) Corrigir erros ortográficos no texto.
- (d) Aumentar o tamanho do corpus de dados.

2. Qual das seguintes etapas faz parte do pré-processamento de texto?

- (a) Tokenização.
- (b) Treinamento do modelo.
- (c) Avaliação do desempenho do modelo.
- (d) Geração de dados sintéticos.

3. O que é tokenização no contexto do pré-processamento de texto?

- (a) O processo de corrigir erros ortográficos em um texto.
- (b) O processo de dividir um texto em palavras, frases ou outros elementos significativos.
- (c) O processo de remover palavras irrelevantes de um texto.
- (d) O processo de transformar texto em vetores numéricos.

4. Qual técnica de pré-processamento reduz palavras às suas formas base ou raízes

- (a) Tokenização.
- (b) Lematização e Stemming.
- (c) Vetorização.
- (d) Filtragem de stopwords.

- 5. Qual biblioteca Python é comumente usada para realizar a tokenização e outras tarefas de pré-processamento de texto?**
- (a) NumPy.
 - (b) Pandas.
 - (c) NLTK.
 - (d) Matplotlib.

Capítulo 4

Vetorização de Texto

Neste capítulo, exploraremos as diferentes técnicas de vetorização de texto, que são essenciais para transformar dados textuais em uma forma numérica que pode ser utilizada por algoritmos de aprendizado de máquina. Abordaremos o *Bag of Words* (BoW), *TF-IDF* e *Word Embeddings*, com exemplos práticos em Python.

4.1 Bag of Words (BoW)

O *Bag of Words* (BoW) é uma das técnicas mais simples de vetorização de texto. Ele representa um texto como um conjunto de palavras, desconsiderando a ordem das palavras, mas mantendo a multiplicidade. Ele simplifica a representação textual para que possa ser utilizada em tarefas de aprendizado de máquina. Imagine um texto como um saco onde jogamos todas as palavras, desconsiderando a ordem em que elas aparecem e contando apenas a frequência com que cada uma delas ocorre.

Como funciona:

1. Tokenização: O texto é dividido em unidades individuais, geralmente palavras.
2. Criação do Vocabulário: É criado um conjunto de todas as palavras únicas presentes em um corpus de textos.

3. Representação Vetorial: Cada documento é representado como um vetor numérico, onde cada posição corresponde a uma palavra do vocabulário e o valor numérico indica a frequência dessa palavra no documento.

Por exemplo, considere dois documentos:

- Documento 1: "O gato subiu na árvore."
- Documento 2: "O cachorro latindo para o gato."

O vocabulário seria: o, gato, subiu, na, árvore, cachorro, latindo, para. A representação vetorial dos documentos poderia ser:

- Documento 1: [1, 2, 1, 1, 1, 0, 0, 0]
- Documento 2: [1, 1, 0, 0, 0, 1, 1, 1]

O Bag-of-Words é fácil de implementar e entender; permite a aplicação de algoritmos de aprendizado de máquina em grandes volumes de texto; além disso, pode ser utilizado em diversas tarefas de PLN, como classificação de textos, clustering e recuperação de informação.

Algumas Limitações é que ele ignora a ordem das palavras e a estrutura gramatical, o que pode levar à perda de significado; além disso existe uma dificuldade em lidar com sinônimos e palavras ambíguas, pois, palavras com significados semelhantes podem ser tratadas como diferentes e palavras com múltiplos significados podem causar ambiguidade.

Algumas aplicações são: a classificação de textos para identificar a categoria de um texto (por exemplo, spam ou não spam); a clusterização, ou seja, agrupar documentos semelhantes; outra é a recuperação de informação para encontrar documentos relevantes em resposta a uma consulta; além disso, a análise de sentimentos para determinar a polaridade de um texto (positivo, negativo ou neutro).

4.1.1 Exemplo em Python: Bag of Words

Código Python

```
1 from sklearn.feature_extraction.text import
    CountVectorizer
2
3 # Exemplo de documentos
4 documentos = [
5     "existe gato laranja",
6     "existe gato preto",
7     "também existe gato branco"
8 ]
9
10 # Criação do modelo BoW
11 vectorizer = CountVectorizer()
12 X = vectorizer.fit_transform(documentos)
13
14 # Exibição da matriz BoW
15 print("Vetor de características:\n", vectorizer.
16       get_feature_names_out())
17 print("Matriz BoW:\n", X.toarray())
```

Saída do console

```
1 Vetor de características:
2 ['branco' 'existe' 'gato' 'laranja' 'preto' 'também']
3 Matriz BoW:
4 [[0 1 1 1 0 0]
5  [0 1 1 0 1 0]
6  [1 1 1 0 0 1]]
```

Este código gera uma matriz BoW, onde cada linha representa um documento e cada coluna representa uma palavra única do vocabulário.

4.2 TF-IDF

O *TF-IDF* é uma técnica de vetorização que leva em consideração a frequência das palavras em um documento em relação a um corpus de documentos. Ele pondera as palavras de acordo com sua importância relativa, diminuindo o peso das palavras comuns e aumentando o peso das palavras menos frequentes. Em termos simples, o TF-IDF nos ajuda a entender quais palavras são mais relevantes e distintivas em um texto específico, comparando-o com outros textos. Utilizando essa técnica, é possível extrair uma métrica estatística para avaliar a importância de uma palavra em um documento em relação a um conjunto de documentos, permitindo assim analisar e entender o conteúdo de grandes volumes de texto.

TF (Term Frequency): Mede a frequência com que uma palavra aparece em um documento. Quanto mais vezes uma palavra aparece, maior é seu TF. IDF (Inverse Document Frequency): Mede a raridade de uma palavra em um conjunto de documentos. Palavras que aparecem em muitos documentos têm um IDF baixo, enquanto palavras que aparecem em poucos documentos têm um IDF alto.

O valor TF-IDF é o produto do TF e do IDF:

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) * \text{IDF}(t)$$

Onde:

- 't' é um termo (palavra)
- 'd' é um documento

Por exemplo, imagine um conjunto de documentos sobre carros. A palavra "carro" provavelmente terá um TF alto em todos os documentos, mas um IDF baixo, pois é muito comum. Já a palavra "supercarro" pode ter um TF baixo em muitos documentos, mas um IDF alto, pois é menos comum. O TF-IDF da palavra "supercarro" será maior do que o da palavra "carro", indicando que "supercarro" é mais distintiva e relevante para identificar documentos sobre carros esportivos.

Alguns exemplos do uso de TF-IDF: recuperação de informação, para ranquear documentos em resposta a uma consulta de pesquisa; mineração

de texto, para identificar tópicos e padrões em grandes conjuntos de documentos. E sistemas de recomendação, para sugerir itens relevantes aos usuários. Ao combinar informações sobre a frequência de palavras em um documento e sua raridade em um conjunto de documentos, o TF-IDF nos permite identificar as palavras mais importantes e relevantes para um determinado contexto.

4.2.1 Exemplo em Python: TF-IDF

Código Python

```
1 from sklearn.feature_extraction.text import
   TfIdfVectorizer
2
3 # Exemplo de documentos
4 documentos = [
5     "existe gato laranja",
6     "existe gato preto",
7     "também existe gato branco"
8 ]
9
10 # Criação do modelo TF-IDF
11 tfidf_vectorizer = TfIdfVectorizer()
12 X_tfidf = tfidf_vectorizer.fit_transform(documentos)
13
14 # Exibição da matriz TF-IDF
15 print("Vetor de características:\n", tfidf_vectorizer.
       get_feature_names_out())
16 print("Matriz TF-IDF:\n", X_tfidf.toarray())
```

Código Python

```
1 Vetor de características:
2 ['branco' 'existe' 'gato' 'laranja' 'preto' 'também']
3 Matriz TF-IDF:
4 [[0.          0.45329466 0.45329466 0.76749457 0.
5          0.          ]
6 [0.          0.45329466 0.45329466 0.
7          0.76749457 0.          ]
8 [0.6088451  0.35959372 0.35959372 0.          0.
9          0.6088451  ]]
```

Este código gera uma matriz TF-IDF, onde cada valor representa a importância de uma palavra em um documento em relação ao corpus.

4.3 Word Embeddings

Os *Word Embeddings* são representações densas de palavras em um espaço vetorial, que capturam as relações semânticas entre as palavras. Técnicas como *Word2Vec*, *GloVe* e *FastText* são amplamente utilizadas para gerar embeddings que refletem o contexto em que as palavras aparecem.

4.3.1 Exemplo em Python: Word2Vec

Código Python

```
1 from gensim.models import Word2Vec
2
3 # Exemplo de corpus
4 corpus = [
5     ["existe", "gato", "laranja"],
6     ["existe", "gato", "preto"],
7     ["também", "existe", "gato", "branco"]
8 ]
9
10 # Treinamento do modelo Word2Vec
11 model = Word2Vec(sentences=corpus, vector_size=100,
12                  window=5, min_count=1, workers=4)
13
14 # Obtenção do embedding para uma palavra
15 print("Embedding para 'gato':\n", model.wv["gato"])
16
17 # Similaridade entre palavras
18 print("Similaridade entre 'gato' e 'laranja':", model.
19       wv.
20       similarity("gato", "laranja"))
```

Código Python

```
1 Embedding para 'gato':
2 [-5.3622725e-04  2.3643136e-04  5.1033497e-03
3    9.0092728e-03
4    ...
5    -8.9173904e-03 -7.0415605e-03  9.0145587e-04
6    6.3925339e-03]
7 Similaridade entre 'gato' e 'laranja': -0.05987629
```

Este exemplo treina um modelo *Word2Vec* e calcula a similaridade entre os vetores de duas palavras.

Em resumo, a vetorização de texto é um passo crucial na preparação de dados para modelos de aprendizado de máquina. Técnicas como BoW, TF-

IDF e Word Embeddings fornecem maneiras eficazes de transformar texto em uma forma que possa ser utilizada para diversas tarefas de processamento de linguagem natural.

Exercícios

Versão on-line destes exercícios

<https://forms.gle/euESDhq7hYWocSGSA>

1. **Qual é a principal característica do modelo Bag of Words (BoW)?**
 - (a) Considera a ordem das palavras no texto.
 - (b) Ignora a ordem das palavras, mas mantém a contagem de frequência das palavras.
 - (c) Gera representações vetoriais densas das palavras.
 - (d) Utiliza embeddings contextuais para representar palavras.
2. **O que o termo TF-IDF representa?**
 - (a) Transform Frequency-Inverse Document Frequency.
 - (b) Term Frequency-Inverse Document Frequency.
 - (c) Term Frequency-Indexed Document Frequency.
 - (d) Transform Frequency-Indexed Document Frequency.
3. **Qual é uma vantagem dos Word Embeddings em relação ao modelo Bag of Words?**
 - (a) Word Embeddings capturam a ordem das palavras no texto.
 - (b) Word Embeddings geram vetores esparsos de alta dimensionalidade.
 - (c) Word Embeddings capturam relações semânticas entre palavras.
 - (d) Word Embeddings ignoram a frequência das palavras no texto.
4. **Qual biblioteca Python tem uma implementação do TF-IDF?**
 - (a) NumPy.
 - (b) Pandas.
 - (c) Scikit-learn.
 - (d) RE.

- 5. Qual técnica é utilizada para criar representações vetoriais densas que capturam o significado semântico das palavras?**
- (a) Bag of Words.
 - (b) TF-IDF.
 - (c) Word Embeddings.
 - (d) N-grams.

Capítulo 5

Modelos Avançados de Vetorização

Neste capítulo, vamos explorar modelos avançados de vetorização de texto, que vão além das técnicas tradicionais como Bag of Words e TF-IDF. Vamos discutir embeddings contextuais, como o BERT e o GPT, e a aplicação de redução de dimensionalidade usando técnicas como PCA e t-SNE. A seguir, forneceremos exemplos práticos em Python para cada um desses métodos.

5.1 Embeddings Contextuais

Embeddings contextuais são vetores que capturam o significado de uma palavra com base no contexto em que ela aparece. Diferente de embeddings como Word2Vec e GloVe, que geram uma única representação para cada palavra, modelos como BERT (Bidirectional Encoder Representations from Transformers) e GPT (Generative Pretrained Transformer) produzem diferentes embeddings para a mesma palavra, dependendo do seu contexto.

5.1.1 Exemplo em Python: Usando BERT para Vetorização

Código Python

```
1 from transformers import BertTokenizer, BertModel
2 import torch
3
4 # Carregando o tokenizer e o modelo BERT
5 tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
6 model = BertModel.from_pretrained("bert-base-uncased")
7
8 # Exemplo de texto
9 texto = "O gato está sentado no tapete."
10
11 # Tokenização do texto e conversão para tensores
12 inputs = tokenizer(texto, return_tensors="pt")
13
14 # Obtenção dos embeddings a partir do modelo BERT
15 with torch.no_grad():
16     outputs = model(**inputs)
17
18 # Extraíndo os embeddings da última camada oculta
19 embeddings = outputs.last_hidden_state
20 print("Embeddings para cada token:", embeddings)
```

Saída do Console

```
1 Embeddings para cada token: tensor([[[[-8.4815e-01,
2     -3.2018e-01, -6.1327e-02, ..., -3.2507e-01,
3     7.6626e-02,  9.1921e-01],
4     [-1.1552e+00, -8.0612e-02, -3.9960e-01, ..., -1.1649e-01,
5     9.8960e-02,  1.0148e+00],
6     ...,
7     [ 6.6603e-01, -7.9354e-02,  1.1458e-02, ..., -4.0039e-02,
8     -7.0081e-01,  2.2797e-02]]]])
```

Este exemplo mostra como usar o BERT para gerar embeddings contextuais. O modelo BERT é capaz de gerar um vetor de embeddings para cada token no texto, considerando o contexto de toda a frase.

5.2 Análise de Sentimento com Embeddings

Modelos de embeddings contextuais, como o BERT, podem ser utilizados em tarefas de análise de sentimento, fornecendo representações mais ricas e precisas das palavras. A seguir, aplicamos BERT em uma tarefa de classificação de sentimento.

5.2.1 Exemplo em Python: Classificação de Sentimento com BERT

Código Python

```
1 from transformers import BertTokenizer,
   BertForSequenceClassification
2 from torch.nn.functional import softmax
3
4 # Carregando o tokenizer e o modelo BERT para
   classificação de sequência
5 tokenizer = BertTokenizer.from_pretrained('bert-base-
   uncased')
6 model = BertForSequenceClassification.from_pretrained(
   'bert-base-uncased', num_labels=2)
7
8 # Exemplo de texto
9 texto = "Eu adoro gatos, eles são fofos"
10
11 # Tokenização do texto e conversão para tensores
12 inputs = tokenizer(texto, return_tensors='pt')
13
14 # Fazendo a previsão
15 with torch.no_grad():
16     outputs = model(**inputs)
17
18 # Aplicando softmax para obter as probabilidades das
   classes
19 probabilidades = softmax(outputs.logits, dim=1)
20 print("Probabilidade de sentimento positivo:",
   probabilidades[0][1].item())
```

Saída do Console

```
1 Probabilidade de sentimento positivo:
   0.6633508801460266
```

Este código ilustra como utilizar o BERT para classificar o sentimento de um texto. Aqui, o modelo é capaz de prever a probabilidade de um

sentimento positivo ou negativo para a frase fornecida.

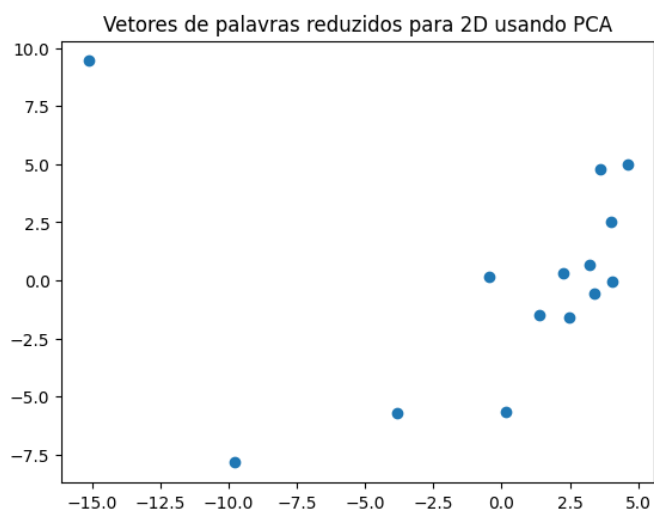
5.3 Redução de Dimensionalidade

Vetores de alta dimensionalidade podem ser difíceis de manipular e visualizar. Técnicas como Análise de Componentes Principais (PCA) e t-SNE (t-Distributed Stochastic Neighbor Embedding) são comumente usadas para reduzir a dimensionalidade dos dados de forma a manter as informações mais importantes.

5.3.1 Exemplo em Python: Redução de Dimensionalidade com PCA

Código Python

```
1 from sklearn.decomposition import PCA
2 import matplotlib.pyplot as plt
3 from transformers import BertTokenizer, BertModel
4 import torch
5
6 # Carregando o tokenizer e o modelo BERT
7 tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
8 model = BertModel.from_pretrained("bert-base-uncased")
9
10 # Exemplo de texto
11 texto = "Eu adoro gatos, eles são fofos"
12
13 # Tokenização do texto e conversão para tensores
14 inputs = tokenizer(texto, return_tensors="pt")
15
16 # Obtenção dos embeddings a partir do modelo BERT
17 with torch.no_grad():
18     outputs = model(**inputs)
19
20 # Extraindo os embeddings da última camada oculta
21 embeddings = outputs.last_hidden_state
22
23 # Exemplo de vetores de alta dimensionalidade (
    usaremos os embeddings do BERT)
24 vetores = embeddings.squeeze().numpy() # Convertendo
    de tensor para numpy array
25
26 # Aplicando PCA para reduzir para 2 dimensões
27 pca = PCA(n_components=2)
28 vetores_reduzidos = pca.fit_transform(vetores)
29
30 # Plotando os vetores em 2D
31 plt.scatter(vetores_reduzidos[:, 0], vetores_reduzidos
   [:, 1])
32 plt.title("Vetores de palavras reduzidos para 2D
    usando PCA")
33 plt.show()
```

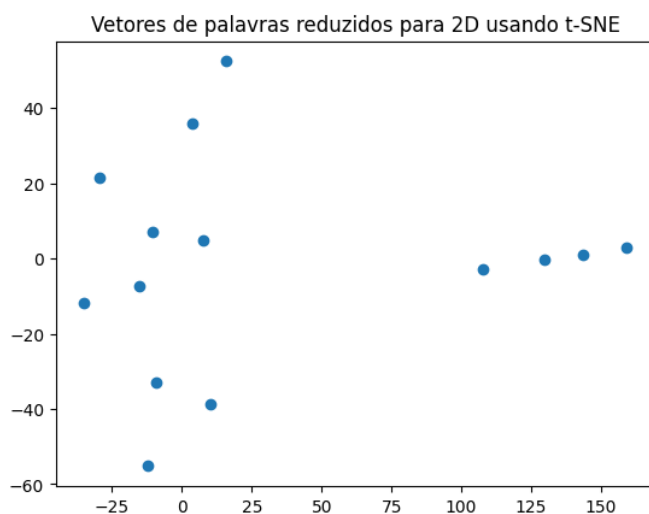


Neste exemplo, utilizamos PCA para reduzir os embeddings gerados pelo BERT para duas dimensões, facilitando a visualização.

5.3.2 Exemplo em Python: Redução de Dimensionalidade com t-SNE

Código Python

```
1 from sklearn.manifold import TSNE
2 import matplotlib.pyplot as plt
3 from transformers import BertTokenizer, BertModel
4 import torch
5
6 # Carregando o tokenizer e o modelo BERT
7 tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
8 model = BertModel.from_pretrained("bert-base-uncased")
9
10 # Exemplo de texto
11 texto = "Eu adoro gatos, eles são fofos"
12
13 # Tokenização do texto e conversão para tensores
14 inputs = tokenizer(texto, return_tensors="pt")
15
16 # Obtenção dos embeddings a partir do modelo BERT
17 with torch.no_grad():
18     outputs = model(**inputs)
19
20 # Extraíndo os embeddings da última camada oculta
21 embeddings = outputs.last_hidden_state
22
23 vetores = embeddings.squeeze().numpy() # Convertendo
    de tensor para numpy array
24
25 # Aplicando t-SNE para reduzir para 2 dimensões
26 tsne = TSNE(n_components=2, perplexity=2)
27 vetores_reduzidos_tsne = tsne.fit_transform(vetores)
28
29 # Plotando os vetores em 2D
30 plt.scatter(vetores_reduzidos_tsne[:, 0],
    vetores_reduzidos_tsne[:, 1])
31 plt.title("Vetores de palavras reduzidos para 2D
    usando t-SNE")
32 plt.show()
```



Aqui, usamos t-SNE para reduzir os embeddings para duas dimensões. O t-SNE é especialmente útil para a visualização de dados em espaços de alta dimensionalidade.

Em resumo, exploramos modelos avançados de vetorização de texto, incluindo embeddings contextuais com BERT e GPT, e técnicas de redução de dimensionalidade como PCA e t-SNE. Esses métodos fornecem ferramentas poderosas para capturar e visualizar informações complexas em dados textuais.

Exercício

Versão on-line destes exercícios

<https://forms.gle/hh4BgCZVhJLMYfiy6>

1. **Qual é a principal vantagem dos embeddings contextuais, como BERT, em comparação com embeddings tradicionais Word2Vec?**
 - (a) Embeddings contextuais capturam o significado das palavras com base em seu contexto específico.
 - (b) Embeddings contextuais são sempre mais rápidos de treinar.
 - (c) Embeddings contextuais geram representações esparsas das palavras.
 - (d) Embeddings contextuais são menos precisos do que os embeddings tradicionais.
2. **O que a técnica de PCA (Análise de Componentes Principais) realiza em dados de alta dimensionalidade?**
 - (a) Aumenta o número de dimensões nos dados.
 - (b) Reduz a dimensionalidade dos dados mantendo a maior variabilidade possível.
 - (c) Gera novas características que não são correlacionadas.
 - (d) Elimina completamente a variabilidade dos dados.
3. **Qual é a função principal da técnica t-SNE?**
 - (a) Agrupar dados de alta dimensionalidade.
 - (b) Visualizar dados de alta dimensionalidade em espaços de menor dimensão.
 - (c) Normalizar dados textuais.
 - (d) Criar novas features a partir dos dados originais.
4. **Por que técnicas de redução de dimensionalidade, como PCA, são úteis em vetorização de texto?**

- (a) Elas aumentam a precisão dos modelos de aprendizado de máquina.
- (b) Elas eliminam a necessidade de embeddings contextuais.
- (c) Elas reduzem a quantidade de dados de entrada para tornar o processamento mais eficiente e visualizável.
- (d) Elas criam novas features para melhorar o desempenho de modelos de deep learning.

5. Em qual cenário a redução de dimensionalidade é particularmente útil?

- (a) Quando se deseja aumentar a complexidade do modelo.
- (b) Quando os dados têm poucas features e baixa variabilidade.
- (c) Quando os dados têm muitas dimensões e se deseja melhorar a visualização ou performance do modelo.
- (d) Quando se deseja eliminar o ruído dos dados, independentemente da dimensionalidade.

Capítulo 6

Aplicações Práticas

Neste capítulo, exploraremos aplicações práticas das técnicas de vetorização de texto que discutimos nos capítulos anteriores. Vamos focar em três áreas principais: classificação de texto, agrupamento de documentos e detecção de tópicos. Cada seção incluirá exemplos práticos em Python.

6.1 Classificação de Texto

A classificação de texto é uma das aplicações mais comuns da vetorização de texto. Milhares de pessoas visitam as mídias sociais para expressar seus sentimentos, principalmente no antigo Twitter, hoje X, [17]. Essas mídias reúnem características que viabilizam a mineração de texto tais como: mensagens textuais, de perfil público e coleta automatizada [10].

A enorme quantidade de dados gerados por essas mídias e as opiniões e sentimentos expressos podem ser analisadas por empresas para diversos fins, sendo possível capturar os sentimentos, opiniões e críticas sobre os produtos em tempo real, proporcionando à empresa uma nova forma de entendimento sobre o comportamento do consumidor [15]. Algumas questões podem ser respondidas através da análise dos milhares de comentários e respostas expressos nessas mídias. Existem evidências que essas postagens são utilizadas na tomada de decisão e que podem influenciar eventos acontecendo em tempo real [3, 20].

As mídias sociais também vêm sendo utilizadas para inferências e previsões em vários setores, por exemplo: já foram encontradas relações entre o sentimento impresso nas postagens dessas mídias e os movimentos do preço do Bitcoin [18, 14]; já foram utilizadas para prever quais seriam as palavras-chave (hashtags) mais comentadas no futuro [6]; para inferir a quantidade de pessoas que irão assistir ao filme [22]; até para correlacionar com os indicados ao Oscar [21]; para influenciar a visão do profissional no ambiente de trabalho, pois, quando os profissionais utilizam as mídias sociais nas relações profissionais e adotam um tipo de comportamentos de gerenciamento de público e de conteúdo, dados apontam uma tendência do profissional ser mais respeitados por seus colegas de trabalho [7].

As mídias sociais também fornecem aos usuários da Internet uma maneira fácil e barata de se envolver em discussões políticas e promover pontos de vista e interesses. Essas postagens também podem auxiliar no processo de decisão de instituições públicas [8]. Alguns candidatos utilizam as mídias sociais para fazer propaganda e para aumentar a intenção de voto dos seus eleitores, além disso, foram encontradas evidências de que essa intenção é influenciada pela opinião do eleitor, confiança e imagem do candidato [2]. O uso da mídia social também é utilizado para divulgar notícias falsas, e existem evidências de que a exposição a essas notícias pode influenciar os resultados das eleições. Nas eleições presidenciais dos EUA em 2016, notícias falsas postadas em mídias sociais foram cruciais na eleição do presidente Trump [1].

Métodos automáticos de extração de sentimentos, também chamado de mineração de opinião, já foram utilizados para avaliar o desempenho das ações na bolsa [5]. O apoio a decisão baseado em sentimentos expresso em mídias sociais tem potencial para beneficiar e diminuir os custos organizacionais, mas ainda é preciso entender como utilizá-la e qual o seu impacto [19].

Usaremos a técnica TF-IDF para vetorizar documentos e, em seguida, aplicaremos um classificador de Naive Bayes para categorizar o texto e analisar o sentimento de algumas postagens, como exemplo.

6.1.1 Exemplo em Python: Classificação de Texto com TF-IDF e Naive Bayes

Código Python

```
1 from sklearn.feature_extraction.text import
    TfidfVectorizer
2 from sklearn.naive_bayes import MultinomialNB
3 from sklearn.pipeline import make_pipeline
4 from sklearn.model_selection import train_test_split
5 from sklearn import metrics
6
7 # Exemplo de corpus e rótulos
8 documentos = [
9     "Este é um ótimo produto",
10    "Muito ruim, não gostei",
11    "Fantástico, recomendo!",
12    "Não vale a pena, péssimo",
13    "Excelente qualidade, compraria de novo",
14    "Horrível, joguei meu dinheiro fora"
15 ]
16 rótulos = ["positivo", "negativo", "positivo", "
    negativo", "positivo", "negativo"]
17
18 # Dividindo os dados em treino e teste
19 X_train, X_test, y_train, y_test = train_test_split(
    documentos, rótulos, test_size=0.33, random_state
    =42)
20
21 # Criando o pipeline
22 modelo = make_pipeline(TfidfVectorizer(),
    MultinomialNB())
23
24 # Treinando o modelo
25 modelo.fit(X_train, y_train)
26
27 # Fazendo previsões
28 predições = modelo.predict(X_test)
29
30 # Avaliando o modelo
31 print(metrics.classification_report(y_test, predições)
    )
```

Código Python

```

1 precision      recall  f1-score   support
2
3 negativo        0.50      1.00      0.67      1
4 positivo        0.00      0.00      0.00      1
5
6 accuracy
7 macro avg       0.25      0.50      0.33      2
8 weighted avg    0.25      0.50      0.33      2

```

Este exemplo demonstra como classificar textos usando TF-IDF e Naive Bayes. Após treinar o modelo, as previsões são feitas no conjunto de teste e uma avaliação do desempenho é realizada.

6.2 Agrupamento de Documentos

O agrupamento de documentos é outra aplicação importante da vetorização de texto, permitindo que documentos similares sejam automaticamente agrupados. Nesta seção, utilizaremos a técnica de *k-means clustering*.

6.2.1 Exemplo em Python: Agrupamento com k-means

Código Python

```
1 from sklearn.feature_extraction.text import
   TfidfVectorizer
2 from sklearn.cluster import KMeans
3
4 # Exemplo de documentos
5 documentos = [
6     "O gato gosta de peixe",
7     "O cachorro gosta de osso",
8     "O peixe não gosta de gato",
9     "O gato e o cachorro são amigos",
10    "O peixe vive na água",
11    "Os cães gostam de ossos"
12 ]
13
14 # Vetorizando os documentos com TF-IDF
15 vectorizer = TfidfVectorizer()
16 X = vectorizer.fit_transform(documentos)
17
18 # Aplicando k-means clustering
19 kmeans = KMeans(n_clusters=2, random_state=42)
20 kmeans.fit(X)
21
22 # Mostrando os rótulos dos clusters
23 print("Rótulos dos Clusters:", kmeans.labels_)
```

Saída do Console

```
1 Rótulos dos Clusters: [0 0 0 0 0 1]
```

Este código realiza o agrupamento dos documentos em dois clusters usando k-means. Cada documento é atribuído a um cluster com base na similaridade de seu conteúdo.

6.3 Detecção de Tópicos

A detecção de tópicos permite identificar os principais assuntos discutidos em um conjunto de documentos. Uma técnica popular para isso é a *Latent Dirichlet Allocation* (LDA), que vamos explorar nesta seção.

6.3.1 Exemplo em Python: Detecção de Tópicos com LDA

Código Python

```
1 from sklearn.decomposition import
    LatentDirichletAllocation
2 from sklearn.feature_extraction.text import
    CountVectorizer
3
4 # Exemplo de documentos
5 documentos = [
6     "O gato gosta de peixe",
7     "O cachorro gosta de osso",
8     "O peixe não gosta de gato",
9     "O gato e o cachorro são amigos",
10    "O peixe vive na água",
11    "Os cães gostam de ossos"
12 ]
13
14 # Vetorizando os documentos com CountVectorizer
15 vectorizer = CountVectorizer()
16 X = vectorizer.fit_transform(documentos)
17
18 # Aplicando LDA
19 lda = LatentDirichletAllocation(n_components=2,
    random_state=42)
20 lda.fit(X)
21
22 # Mostrando os tópicos
23 tópicos = lda.components_
24 nomes_palavras = vectorizer.get_feature_names_out()
25
26 for idx, tópico in enumerate(tópicos):
27     print(f"Tópico {idx + 1}:")
28     palavras = [nomes_palavras[i] for i in tópico.
        argsort()[: -5 - 1 : -1]]
29     print(" ".join(palavras))
```


Saída do Console

```
1 Tópico 1:  
2 gato peixe gosta de cachorro  
3 Tópico 2:  
4 de ossos os gostam cães
```

Este exemplo demonstra como usar LDA para identificar tópicos em um conjunto de documentos. As palavras mais representativas para cada tópico são exibidas.

Em resumo, as técnicas de vetorização de texto discutidas anteriormente têm uma ampla gama de aplicações práticas, desde a classificação e agrupamento de documentos até a detecção de tópicos. Essas aplicações ilustram a importância de transformar texto em representações numéricas para a análise eficiente de dados textuais.

Exercícios

Versão on-line destes exercícios

<https://forms.gle/FopDZjLXRvixmGW46>

1. **Qual é o principal objetivo da classificação de texto?**
 - (a) Gerar texto automaticamente.
 - (b) Classificar textos em categorias predefinidas.
 - (c) Traduzir textos entre diferentes idiomas.
 - (d) Aumentar o tamanho do corpus de dados textuais.
2. **Qual técnica de vetorização é comumente usada em combinação para classificação de texto?**
 - (a) Dinâmica de Sistemas.
 - (b) TF-IDF ou BoW.
 - (c) Redes de Petri.
 - (d) Automação.
3. **O que o algoritmo de k-means clustering faz em um conjunto de documentos?**
 - (a) Agrupa documentos semelhantes em clusters.
 - (b) Classifica documentos em categorias predefinidas.
 - (c) Reduz a dimensionalidade dos vetores de documentos.
 - (d) Gera novos documentos a partir de um conjunto de treinamento.
4. **Qual técnica é utilizada para identificar os principais tópicos em um conjunto de documentos?**
 - (a) Classificação de texto.
 - (b) Redução de dimensionalidade.
 - (c) Modelagem de tópicos com LDA.
 - (d) Tokenização.

5. **Em um modelo de classificação de sentimentos, qual é o propósito de utilizar TF-IDF na vetorização de texto?**
- (a) Melhorar a visualização dos dados.
 - (b) Reduzir o número de palavras no texto.
 - (c) Ponderar a importância das palavras com base na frequência no documento e no corpus.
 - (d) Gerar novos tópicos a partir dos documentos.

Capítulo 7

Estudo de Caso

Neste capítulo, exploraremos alguns estudos de caso que demonstram a aplicação prática das técnicas de vetorização de texto discutidas anteriormente. Cada caso de estudo será acompanhado por exemplos em Python para ilustrar como essas técnicas podem ser implementadas em situações do mundo real.

7.1 Análise de Reviews de Produtos

Uma aplicação comum da vetorização de texto é a análise de opiniões e reviews de produtos. Este caso de estudo examina como podemos utilizar TF-IDF e modelos de classificação para analisar reviews de produtos e determinar a polaridade dos sentimentos expressos.

7.1.1 Exemplo em Python: Análise de Sentimentos em Reviews de Produtos

Código Python

```
1 from sklearn.feature_extraction.text import
    TfidfVectorizer
2 from sklearn.naive_bayes import MultinomialNB
3 from sklearn.pipeline import make_pipeline
4 from sklearn.model_selection import train_test_split
5 from sklearn import metrics
6
7 # Exemplo de reviews de produtos e rótulos
8 reviews = [
9     "Este produto é excelente, superou minhas
        expectativas!",
10    "Horrível, não funcionou como esperado",
11    "Muito bom, compraria novamente",
12    "Péssimo, nunca mais compro deste vendedor",
13    "Produto de ótima qualidade",
14    "Não gostei, material de baixa qualidade"
15 ]
16 rótulos = ['positivo', 'negativo', 'positivo', '
    negativo', 'positivo', 'negativo']
17
18 # Dividindo os dados em treino e teste
19 X_train, X_test, y_train, y_test = train_test_split(
    reviews, rótulos, test_size=0.33, random_state=42)
20
21 # Criando o pipeline TF-IDF + Naive Bayes
22 modelo = make_pipeline(TfidfVectorizer(),
    MultinomialNB())
23
24 # Treinando o modelo
25 modelo.fit(X_train, y_train)
26
27 # Fazendo previsões
28 predições = modelo.predict(X_test)
29
30 # Avaliando o modelo
31 print(metrics.classification_report(y_test, predições)
    )
```

Saída do Console					
1		precision	recall	f1-score	support
2					
3	negativo	1.00	1.00	1.00	1
4	positivo	1.00	1.00	1.00	1
5					
6	accuracy			1.00	2
7	macro avg	1.00	1.00	1.00	2
8	weighted avg	1.00	1.00	1.00	2

Este exemplo mostra como aplicar TF-IDF e um classificador de Naive Bayes para analisar reviews de produtos, determinando se o sentimento expresso é positivo ou negativo.

7.2 Processamento de Tweets

Analisar e processar tweets é uma tarefa desafiadora devido à natureza informal e limitada dos textos. Este caso de estudo se concentra na análise de sentimentos e na classificação de tweets, utilizando técnicas de vetorização e aprendizado de máquina.

7.2.1 Exemplo em Python: Análise de Sentimentos em Tweets

Código Python

```
1 import pandas as pd
2 from sklearn.feature_extraction.text import
    TfidfVectorizer
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.model_selection import train_test_split
5 from sklearn import metrics
6
7 # Exemplo de tweets e rótulos
8 tweets = [
9     "Adoro usar o novo recurso, muito útil!",
10    "Não gosto desta atualização, ficou pior",
11    "Essa nova versão está incrível",
12    "Detestei o que fizeram no app",
13    "Excelente trabalho, equipe!",
14    "Uma das piores atualizações até agora"
15 ]
16 rótulos = ['positivo', 'negativo', 'positivo', '
    negativo', 'positivo', 'negativo']
17
18 # Dividindo os dados em treino e teste
19 X_train, X_test, y_train, y_test = train_test_split(
    tweets, rótulos, test_size=0.33, random_state=42)
20
21 # Criando o pipeline TF-IDF + Regressão Logística
22 modelo = make_pipeline(TfidfVectorizer(),
    LogisticRegression())
23
24 # Treinando o modelo
25 modelo.fit(X_train, y_train)
26
27 # Fazendo previsões
28 predições = modelo.predict(X_test)
29
30 # Avaliando o modelo
31 print(metrics.classification_report(y_test, predições)
    )
```

Saída do Console					
1		precision	recall	f1-score	support
2					
3	negativo	0.00	0.00	0.00	1
4	positivo	0.50	1.00	0.67	1
5					
6	accuracy			0.50	2
7	macro avg	0.25	0.50	0.33	2
8	weighted avg	0.25	0.50	0.33	2

Este exemplo ilustra como podemos utilizar TF-IDF e Regressão Logística para analisar tweets, classificando-os como positivos ou negativos com base no conteúdo.

7.3 Análise de Notícias

A análise de notícias é fundamental para entender as tendências atuais e os tópicos de interesse público. Neste caso de estudo, utilizaremos técnicas de vetorização e modelagem de tópicos para identificar os principais temas abordados em um conjunto de artigos de notícias.

7.3.1 Exemplo em Python: Modelagem de Tópicos em Notícias

Código Python

```
1 from sklearn.decomposition import
    LatentDirichletAllocation
2 from sklearn.feature_extraction.text import
    CountVectorizer
3
4 # Exemplo de notícias
5 noticias = [
6     "O governo aprova novas reformas econômicas",
7     "A tecnologia 5G está mudando o cenário global",
8     "Novas descobertas na área da saúde são",
9     "promissoras",
10    "A economia global enfrenta novos desafios",
11    "Tecnologias emergentes como IA estão em alta",
12    "Novas políticas ambientais são implementadas"
13 ]
14 # Vetorizando as notícias com CountVectorizer
15 vectorizer = CountVectorizer()
16 X = vectorizer.fit_transform(noticias)
17
18 # Aplicando LDA para modelagem de tópicos
19 lda = LatentDirichletAllocation(n_components=2,
    random_state=42)
20 lda.fit(X)
21
22 # Mostrando os tópicos
23 tópicos = lda.components_
24 nomes_palavras = vectorizer.get_feature_names_out()
25
26 for idx, tópico in enumerate(tópicos):
27     print(f"Tópico {idx + 1}:")
28     palavras = [nomes_palavras[i] for i in tópico.
        argsort()[::-5 - 1:-1]]
29     print(" ".join(palavras))
```

Saída do Console

```
1 Tópico 1:  
2 global estão alta como em  
3 Tópico 2:  
4 novas são promissoras área descobertas
```

Neste exemplo, aplicamos LDA para identificar tópicos em um conjunto de artigos de notícias, mostrando as palavras mais representativas para cada tópico.

Em resumo, os casos de estudo apresentados neste capítulo demonstram como as técnicas de vetorização de texto podem ser aplicadas a problemas reais, como análise de sentimentos em reviews e tweets, e modelagem de tópicos em notícias. Essas técnicas são ferramentas poderosas para extrair informações valiosas de grandes volumes de dados textuais.

Exercício

Versão on-line destes exercícios

<https://forms.gle/uuFPqkeQFtPREuQy5>

1. **Qual é o principal objetivo da análise de sentimentos?**
 - (a) Aumentar o número de visualizações.
 - (b) Identificar e classificar sentimentos expressos no texto.
 - (c) Melhorar a qualidade dos produtos.
 - (d) Reduzir o tempo de leitura dos reviews.
2. **Qual é o principal desafio ao realizar análise de sentimentos?**
 - (a) A falta de técnicas de vetorização adequadas.
 - (b) O número limitado de palavras únicas nos reviews.
 - (c) A ambiguidade no sentimento expresso em algumas frases.
 - (d) A necessidade de grandes volumes de dados de treinamento.
3. **Qual é o principal desafio ao processar tweets em comparação com outros tipos de textos?**
 - (a) Tweets são sempre escritos de forma gramaticalmente correta.
 - (b) Tweets têm comprimento fixo e linguagem informal, o que pode dificultar a análise.
 - (c) Tweets não contêm informações úteis para análise de sentimentos.
 - (d) Tweets são muito longos e detalhados, o que dificulta a análise.
4. **Qual técnica é utilizada para identificar os principais temas em um conjunto de artigos de notícias?**
 - (a) Classificação de Sentimentos.
 - (b) Análise de Regressão.
 - (c) Modelagem de Tópicos com LDA.

(d) Clustering com k-means.

5. Qual é uma aplicação comum da modelagem de tópicos em notícias?

- (a) Traduzir notícias para diferentes idiomas.
- (b) Identificar e agrupar artigos com temas semelhantes.
- (c) Prever o futuro de tendências de notícias.
- (d) Criar resumos automáticos de notícias.

Próximos Passos

8.1 Futuras Direções no Campo de Vetorização de Texto

O campo da vetorização de texto e PLN está em constante evolução. A seguir, algumas direções futuras e tendências que merecem atenção:

8.1.1 Modelos de Linguagem de Grande Escala

Modelos de linguagem de grande escala, como GPT-3 e modelos de nova geração, estão redefinindo o que é possível em PLN. Eles não só geram texto de alta qualidade, como também podem ser usados para uma variedade de tarefas de NLP, incluindo tradução automática, resumo de textos, e geração de código.

8.1.2 Multimodalidade

A integração de dados textuais com outros tipos de dados, como imagens e áudio, está ganhando força. Modelos multimodais, como CLIP (Contrastive Language-Image Pretraining), permitem o entendimento e a geração de conteúdo que combina múltiplas formas de mídia.

8.1.3 Vetorização de Texto em Tempo Real

Com a necessidade crescente de processamento em tempo real, técnicas para vetorização de texto que possam ser aplicadas instantaneamente a fluxos de dados contínuos estão se tornando cada vez mais relevantes. Aplicações como chatbots, assistentes virtuais e sistemas de recomendação em tempo real se beneficiam dessas abordagens.

8.2 Leitura Recomendada

Para aqueles que desejam se aprofundar ainda mais no assunto, aqui estão algumas leituras recomendadas que cobrem uma gama de tópicos desde os fundamentos até as últimas inovações em vetorização de texto e PLN.

- Speech and Language Processing por Daniel Jurafsky e James H. Martin [12].
- Deep Learning por Ian Goodfellow, Yoshua Bengio e Aaron Courville [9].
- Introduction to Information Retrieval por Christopher D. Manning, Prabhakar Raghavan e Hinrich Schütze [13].
- Natural Language Processing with Python por Steven Bird, Ewan Klein e Edward Loper [4].

8.3 Próximos Passos

Com o conhecimento adquirido neste livro, você está agora equipado para aplicar técnicas de vetorização de texto em projetos de NLP. Alguns próximos passos sugeridos incluem:

- Desenvolver Aplicações Reais: Aplique as técnicas aprendidas para resolver problemas reais, como análise de sentimentos em redes sociais ou categorização automática de e-mails.

- Explorar Novos Modelos: Experimente modelos de linguagem mais recentes, como GPT-3 ou T5, para ver como eles se comparam aos métodos abordados aqui.
- Participar de Comunidades: Engaje-se com a comunidade de PLN e machine learning através de fóruns online, workshops, e competições de codificação, como as organizadas pelo Kaggle.

Esperamos que este livro tenha servido como uma base para seu crescimento contínuo no campo de processamento de linguagem natural e vetorização de texto.

Gabarito das questões

Gabarito das questões do exercício:

Cap 1: 1 B, 2 C, 3 B, 4 B, 5 C

Cap 2: 1 B, 2 C, 3 C, 4 B, 5 C

Cap 3: 1 B, 2 A, 3 B, 4 B, 5 C

Cap 4: 1 B, 2 B, 3 C, 4 C, 5 C

Cap 5: 1 A, 2 B, 3 B, 4 C, 5 C

Cap 6: 1 B, 2 B, 3 A, 4 C, 5 C

Cap 7: 1 B, 2 C, 3 B, 4 C, 5 B

Sobre os autores

GISELDO NEO é Professor de Informática no Instituto Federal de Alagoas (IFAL) campus Viçosa, Doutorando em Ciência da Computação na Universidade Federal de Campina Grande, Mestre em Modelagem Computacional do Conhecimento na Universidade Federal de Alagoas, Mestre em Contabilidade (FUCAPE). Possui MBA em Gestão e Estratégia Empresarial (ESTÁCIO), Especialização em Arquitetura e Engenharia de Software (ESTÁCIO), MBA em Gestão de Projetos (UNINTER). Graduação em Análise e Desenvolvimento de Sistemas (ESTÁCIO), Graduação em Processos Gerenciais (UNINTER) e Técnico de Informática (IFS).

ALANA NEO é Professora de Informática no Instituto Federal do Mato Grosso do Sul (IFMS) e atualmente desenvolve pesquisas na área de Informática na Educação. Doutoranda em Ciência da Computação na Universidade Federal de Campina Grande (UFCG), Mestra em Modelagem Computacional do Conhecimento na Universidade Federal de Alagoas (UFAL), Especialista em Estratégias Didáticas para a Educação Básica com Uso de TIC na Universidade Federal de Alagoas (UFAL), Especialista em Desenvolvimento de Software, Especialista em Segurança da Informação, Graduada em Análise e Desenvolvimento de Sistemas e Bacharel em Sistemas de Informação pela Universidade Estácio de Sá e Licenciatura em Computação pelo Claretiano Centro Universitário.

Referências Bibliográficas

- [1] Hunt Allcott and Matthew Gentzkow. Social media and fake news in the 2016 election. *Journal of Economic Perspectives*, 31(2):211–236, 2017.
- [2] Cláudio Márcio Almeida and Marcelo de Oliveira. *O Impacto da Mídia Social na Intenção de Voto do Eleitor*. Independently Published, Vitória, 1 edition, 2018.
- [3] Baldykowski, Ariane Lao Saulo Alves de Brito, Sandro A. Miczevski, and Thiago H. Silva. Cheers to untappd! preferences for beer reflect cultural differences around the world. *Americas Conference on Information Systems 2018: Digital Disruption, AMCIS 2018*, (July), 2018.
- [4] Steven Bird, Ewan Klein, and Edward Loper. *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*. O’Reilly Media, Inc., 2009.
- [5] Hsinchun Chen and David Zimbra. AI and opinion mining. *IEEE Intelligent Systems*, 25(4):72–79, 2010.
- [6] Anubrata Das, Moumita Roy, Soumi Dutta, Saptarshi Ghosh, and Asit Kumar Das. Predicting Trends in the Twitter Social Network: A Machine Learning Approach. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 8947:432–444, 2015.

- [7] Amanda Soares Zambelli Ferretti and Bruno Felix Von Borell de Araujo. Comportamentos nas Redes Sociais Online e seus Impactos nas Relações Profissionais. *Revista Administração em Diálogo - RAD*, 19(2):91, 2017.
- [8] Elena Georgiadou, Spyros Angelopoulos, and Helen Drake. Big data analytics and international negotiations: Sentiment analysis of Brexit negotiating outcomes. *International Journal of Information Management*, 51(October 2019):102048, 2020.
- [9] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- [10] Hootsuite. Most popular social networks worldwide as of January 2018, ranked by number of active users (in millions), 2018.
- [11] Daniel Jurafsky and James H Martin. *Speech and Language Processing*. Prentice Hall, 2000.
- [12] Daniel Jurafsky and James H Martin. *Speech and Language Processing*. Prentice Hall, 2009.
- [13] Christopher D Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [14] Satoshi Nakamoto. Bitcoin: A Peer-to-Peer Electronic Cash System. 2008.
- [15] Aldomar Nascimento Junior. *Com a boca no twitter: a cocriação e a colaboração impactando no sucesso do marketing boca a boca on-line*. PhD thesis, 2012.
- [16] Alexander Pak and Patrick Paroubek. Twitter as a corpus for sentiment analysis and opinion mining. In *Proceedings of the Seventh Conference on International Language Resources and Evaluation (LREC'10)*, pages 1320–1326. European Language Resources Association (ELRA), 2010.

- [17] Raquel Recuero and Gabriela Zago. Em busca das redes que importam: redes sociais e capital social no Twitter. pages 81–94, 2016.
- [18] Weuler Borges Santos Santos. Prevendo o preço do Bitcoin com Redes Neurais usando dados do Twitter e de Mercado. 2019.
- [19] Gregory D. Saxton and Chao Guo. Social media capital: Conceptualizing the nature, acquisition, and expenditure of social media-based organizational resources. *International Journal of Accounting Information Systems*, 36:100443, 2020.
- [20] David A. Shamma, Lyndon Kennedy, and Elizabeth F. Churchill. Tweet the debates: Understanding community annotation of uncollected sources. *1st ACM SIGMM International Workshop on Social Media, WSM'09, Co-located with the 2009 ACM International Conference on Multimedia, MM'09*, pages 3–10, 2009.
- [21] Igor Tannús Corrêa. Análise dos sentimentos expressos na rede social Twitter em relação aos filmes indicados ao Oscar 2017. 2017.
- [22] Diogo Teixeira and Isabel Azevedo. Análise de opiniões expressas nas redes sociais. *RISTI - Revista Iberica de Sistemas e Tecnologias de Informacao*, (8):53–65, 2011.