



数据分析：文本分析与自然语言处理

文本分析基础

1. 文本数据的预处理

文本预处理是文本分析与自然语言处理(NLP)中至关重要的第一步，它包括清洗、规范化和转换文本数据，以准备用于后续的分析 and 建模。预处理的目的是去除文本中的噪声，如HTML标签、标点符号、数字、特殊字符等，同时将文本转换为一致的格式，如统一大小写、去除多余的空格等。

1.1 示例代码

```
import re
import string

def clean_text(text):
    # 转换为小写
    text = text.lower()
    # 去除数字
    text = re.sub(r'\d+', '', text)
    # 去除标点符号
    text = text.translate(str.maketrans('', '', string.punctuation))
    # 去除多余的空格
    text = re.sub(r'\s+', ' ', text)
    return text

# 示例文本
sample_text = "Hello, World! This is a sample text with numbers 123
and punctuation!!!"
# 清洗文本
cleaned_text = clean_text(sample_text)
print(cleaned_text)
```

1.2 代码解释

上述代码定义了一个`clean_text`函数，用于执行文本清洗任务。它首先将文本转换为小写，然后使用正则表达式去除所有数字和标点符号，最后去除多余的空格。`sample_text`是一个包含大小写字母、数字和标点的示例文本，通过调用`clean_text`函数，可以看到清洗后的结果。

2. 分词与词干提取

分词是将文本分割成单词或短语的过程，而词干提取则是将单词转换为其基本形式或词干，以减少词汇的多样性并提高分析的效率。在英语中，词干提取通常涉及去除词缀，如“running”转换为“run”。

2.1 示例代码

```
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer

# 初始化词干提取器
stemmer = PorterStemmer()

def stem_words(words):
    # 对每个单词进行词干提取
    stemmed_words = [stemmer.stem(word) for word in words]
    return stemmed_words

# 示例文本
sample_text = "I am running and jumping and playing."
# 分词
words = word_tokenize(sample_text)
# 词干提取
stemmed_words = stem_words(words)
print(stemmed_words)
```

2.2 代码解释

这段代码使用了NLTK库中的`word_tokenize`函数进行分词，然后使用`PorterStemmer`进行词干提取。`stem_words`函数接收一个单词列表，对列表中的每个单词应用词干提取器，最后返回词干化的单词列表。

3. 停用词的去除

停用词是指在文本中频繁出现但对分析意义不大的词汇，如“the”、“is”、“at”等。去除停用词可以减少数据的维度，同时提高模型的性能。

3.1 示例代码

```
from nltk.corpus import stopwords

def remove_stopwords(words):
    # 英语停用词列表
```

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：<https://d.book118.com/396031001114010201>