囯

### D09 基礎語言模型:N-Gram



#### 重要知識點



- 透過實作更加了解N-gram模型
- 使用NLTK建構N-gram模型

### N-Gram語言模型

- 字當作條件機率計算的先驗條件,形成長度是N的字詞片段序列。每個字詞片段及稱為gram。
  - 在了解N-Gram語言模型的原理基礎後,現在來看該如何使用Python來實踐N-Gram模型。

• N-Gram語言模型是基於統計的語言模型算法,基於馬可夫假設將文本中的內容取最靠近的N個

### 建立詞頻字典

N-Gram語言模型是基於詞頻來計算機率,因此須先將文本資料轉化為詞頻字典。

- 準備文本資料(這裡將文本資料儲存在words變數中)
- 這裡的preprocess為將文本資料做前處理(詳細可見附件.ipynb)

### words 會是這樣的資料型態 ['the', 'project', 'gutenberg', 'of', 'pride', ...]

words = preprocess(corpus)

建立詞頻字典(bigram)

1 # unigram 2 unigram\_frequecy = dict() 4 for word in words: 5 unigram\_frequecy[word] = unigram\_frequecy.get(word, 0) + 1

7 # 根據詞類排序, 並轉換為(word, count)格式 8 unigram\_frequecy = sorted(unigram\_frequecy.items(), key=lambda word\_count: word\_count[1], reverse=True) 10 # 查看詞類前10的字詞 11 print(unigram\_frequecy[:10]) [('the', 4507), ('to', 4243), ('of', 3730), ('and', 3658), ('her', 2225), ('i', 2070), ('a', 2012), ('in', 1937), ('w as', 1847), ('she', 1710)] 1 # unigram

unigram\_frequecy = dict() 4 for word in words: unigram\_frequecy[word] = unigram\_frequecy.get(word, 0) + 1 7 # 根據詞頻排序, 並轉換為(word, count)格式 8 unigram\_frequecy = sorted(unigram\_frequecy.items(), key=lambda word\_count: word\_count[1], reverse= 10 # 查看詞頻前10的字詞 11 print(unigram\_frequecy[:10])

### 搭建N-Gram模型

```
建立Bigram模型
```

```
def do_bigram_prediction(bigram_freq, start_word, num_words):
#定義起始字
       pred_words = [start_word]
       word = start_word
5 for i in range(num_words):
         # 找尋下一個字
           word = next((word_pairs[1] for (word_pairs, count) in bigram_freq if word_pairs[0] == word), None)
           if not word:
10
              break
           else:
              pred_words.append(word)
13 return pred_words
```

def do\_bigram\_prediction(bigram\_freq, start\_word, num\_words): #定義起始字 pred\_words = [start\_word] word = start\_word for i in range(num\_words): # 找尋下一個字 word = next((word\_pairs[1] for (word\_pairs, count) in bigram\_freq if word\_pairs[0] == word if not word: 10 break 11 else: pred words.append(word)

## 使用N-Gram模型進行預測

#### 使用建立好的Bigram模型進行預測 import random

4 start\_word = random.choice(words) 5 print(f'初始字: {start\_word}') 7 # 使用選取的起始字預測接下來的字詞(共10個字) pred\_words = do\_bigram\_prediction(bigram\_frequency, start\_word, 10) 11 print(f"預測句子: {' '.join(pred\_words)}") 預測句子: of the whole party were to be so much as to

1 import random 3 # 隨機選取起始字 4 start\_word = random.choice(words) print(f'初始字: {start\_word}') 7 # 使用選取的起始字預測接下來的字詞(共10個字) pred\_words = do\_bigram\_prediction(bigram\_frequency, start\_word, 10) 10 # 顯示預測結果 11 print(f"預測句子: {' '.join(pred\_words)}")

## 使用NLTK建構N-Gram模型

NLTK(Natural Language Toolkit)使一個python的自然語言處理套件,內含許多應用, 包括了N-Gram模型

### 安裝NLTK套件 pip install nltk

使用NLTK API搭建N-Gram import collections

#### from nltk import ngrams 4 #使用NLTK API搭建Bigram 7 #使用collectins套件計算詞類

5 bigram\_frequency = ngrams(words, n=2) 8 bigram\_frequency = collections.Counter(bigram\_frequency) 9 print(bigram\_frequency) Counter({('of', 'the'): 491, ('to', 'be'): 445, ('in', 'the'): 397, ('i', 'am'): 303, ('mr', 'darcy'): 273, ('to', 'the'): 268, ('of', 'her'): 261, ('it', 'was'): 251, ('of', 'his'): 235, ('she', 'was'): 212, ('she', 'had'): 205, ('ha d', 'been'): 200, ('it', 'is'): 194, ('i', 'have'): 188, ('to', 'her'): 179, ('that', 'he'): 177, ('could', 'not'): 1 67, ('he', 'had'): 166, ('and', 'the'): 165, ('for', 'the'): 163, ('he', 'was'): 160, ('on', 'the'): 154, ('mrs', 'be nnet'): 153, ('mr', 'collins'): 150, ('such', 'a'): 146, ('in', 'a'): 138, ('do', 'not'): 135, ('have', 'been'): 135, ('with', 'the'): 134, ('they', 'were'): 133, ('did', 'not'): 132, ('she', 'could'): 132, ('that', 'she'): 132, ('an d', 'i'): 124, ('was', 'not'): 124, ('my', 'dear'): 120, ('of', 'a'): 116, ('by', 'the'): 116, ('lady', 'catherine'): 116, ('mr', 'bingley'): 115, ('all', 'the'): 115, ('that', 'i'): 114, ('and', 'she'): 112, ('as', 'to'): 105, ('you', 'are'): 103, ('at', 'the'): 103, ('but', 'i'): 102, ('from', 'the'): 97, ('in', 'her'): 95, ('as', 'she'): 95, ('would', 'be'): 93, ('not', 'be'): 93, ('i', 'do'): 92, ('to', 'his'): 92, ('with', 'a'): 92, ('will', 'be'): 90, ('of', 'it'): 90, ('her', 'sister'): 90, ('mr', 'bennet'): 89, ('for', 'her'): 89, ('and', 'her'): 88, ('there', 'was'): 88,

1 import collections 2 from nltk import ngrams 4 #使用NLTK API搭建Bigram bigram\_frequency = ngrams(words, n=2) #使用collectins套件計算詞頻 bigram\_frequency = collections.Counter(bigram\_frequency) print(bigram\_frequency)

# 詳細使用操作

請參照使用基礎語言模型\_Ngram.ipynb檔進行更詳細的使用操作

### 知識點回顧 在這章節我們學習到了

網站:<u>取的N-Gram</u>

使用NLTK建構N-Gram

• 實作與使用N-gram模型 • 使用NLTK建構N-gram模型

延伸閱讀

Copy and Edit 111 Version 1 of 1 The order words are used in is important The order that words are used in text is not random. In English, for example, you can say "the red Notebook apple" but not "apple red the". The relationships between words in text is very complex. So complex, in fact, that there's an entire field of linguistics devoted to it: syntax. (You can learn more about syntax here, if you're interested.) Execution Info However, there is a relatively simply way of capturing some of these relationships between words, Log originally proposed by Warren Weaver. You can capture quite a bit of information by just looking at Comments (9) which words tend to show up next to each other more often. What are ngrams? The general idea is that you can look at each pair (or triple, set of four, etc.) of words that occur next to each other. In a sufficently-large corpus, you're likely to see "the red" and "red apple" several times, but less likely to see "apple red" and "red the". This is useful to know if, for example, you're trying to figure out what someone is more likely to say to help decide between possible output for an automatic speech recognition system. These co-occuring words are known as "n-grams", where "n" is a number saying how long a string of words you considered. (Unigrams are single words, bigrams are two words, trigrams are three words, 4-grams are four words, 5-grams are five words, etc.) Learning goal: In this tutorial, you'll learn how to find all the n-grams in a corpus & count how often each is used.

### 網站:<u>使用Python搭建N-Gram模型應用</u> 利用Python搭建得N-Gram模型進行新聞文本預測應用

用Python 實現n-gram 語言模型進行新聞文本內容預測 Jed · 2019-12-01 · 3評論 · 1593閱讀 概述

### 語言模型與新聞的內容預測, 分為4個階段: 1. 從主流新聞網站科技頻道抓取新聞內容;

- 2. 數據預處理; 3. 實現n-gram 語言模型; 4. 在測試集上檢驗模型。
- 完整代碼: https://github.com/xJed/ngram-text-prediction

網站:<u>NLTK使用教學手冊</u> NLTK各項應用教學

# NLTK (Natural Language Toolkit) Tutorial in Python

Processing? Natural Language Processing is manipulation or

understanding text or speech by any software or machine. An analogy is that humans interact, understand each other views, and respond with the

What is Natural Language

appropriate answer. In NLP, this interaction, understanding, the response is made by a computer instead of a human. NLTK stands for Natural Language Toolkit. This toolkit is one of the most powerful NLP libraries which contains packages to make machines understand human language and reply to it with an

Home Testing SAP Web Must Learn! Big Data Live Projects

◀ Trending

Course

DevOps

 DBMS Al

Jenkins

What is NLTK?

appropriate response. Tokenization, Stemming, Lemmatization, Punctuation, Character count, word count are some of these packages which will be discussed in this tutorial. Here is what we cover in the Course