

Modul 2 Praktikum Natural Language Processing Preprocessing Data Teks menggunakan NLTK

Deskripsi Singkat

Modul ini membahas tentang preprocessing data teks menggunakan Natural Language Toolkit (NLTK). Dalam modul ini, peserta akan mempelajari cara menggunakan regex untuk parsing teks, teknik-teknik penanganan string, dan berbagai metode untuk mengeksplorasi serta memproses data teks sebelum digunakan dalam analisis lebih lanjut.

Tujuan

1. Memahami konsep dasar preprocessing data teks.
2. Menggunakan regex untuk parsing teks.
3. Menerapkan teknik-teknik penanganan string.
4. Mengeksplorasi dan memproses data teks menggunakan NLTK.
5. Menyiapkan data teks untuk analisis lebih lanjut.

Materi 1 - Parsing teks menggunakan regex

RegEx (Regular Expressions) yaitu sekumpulan karakter yang membentuk pola pencarian. Regex sangat berguna dalam berbagai pemrosesan teks, seperti validasi input, pencarian dan penggantian teks, serta ekstraksi informasi dari teks.

Regex	Deskripsi
[ab]	Mencocokkan karakter a dan b
^ab	Mencocokkan karakter dari a sampai d
[a-z]	Mencocokkan karakter dari a sampai z
[A-Z]	Mencocokkan karakter dari A sampai Z
[0-9]	Mencocokkan karakter dari 0 sampai 9
\s	Mencocokkan spasi
\S	Mencocokkan karakter bukan spasi
\d	Mencocokkan digit
\D	Mencocokkan karakter bukan digit
\w	Mencocokkan kata
\W	Mencocokkan selain kata
+	Mencocokkan satu atau lebih karakter
*	Mencocokkan nol atau lebih karakter
?	Mencocokkan nol atau satu karakter

Untuk simbol-simbol lainnya, silahkan kunjungi
https://www.w3schools.com/python/python_regex.asp.

Mengambil email dari sebuah text

```
1 import re
2
3 doc = "Informasi lebih lanjut hubungi: xyz@abc.com, pqr@mno.com"
4
5 addresses = re.findall(r'[\w\.-]+@[\w\.-]+', doc)
6 for address in addresses:
7     print(address)
8
9 # Output:
10 xyz@abc.com
11 pqr@mno.com
```

Berikut adalah penjelasan dari regex yang digunakan:

- `[\w\.-]+` : Mencocokkan satu atau lebih karakter alfanumerik, titik, atau strip.
- `@` : Mencocokkan karakter @.
- `[\w\.-]+` : Mencocokkan satu atau lebih karakter alfanumerik, titik, atau strip.

Merubah email dari sebuah text

```
1 import re
2
3 doc = "Informasi lebih lanjut hubungi: xyz@abc.com"
4
5 email_baru = re.sub(r'([\w\.-]+)@([\w\.-]+)', r'pqr@mno.com', doc)
6 print(email_baru)
7
8 # Output
9 Informasi lebih lanjut hubungi: pqr@mno.com
```

Function `re.sub()` digunakan untuk mengganti teks yang cocok dengan pola tertentu dengan teks lain.

Melakukan tokenisasi teks menggunakan regex

```
1 import re
2
3 doc = "Saya sedang belajar NLP."
4
5 token = re.split(r'\s+', doc)
6
7 # Output:
8 ['Saya', 'sedang', 'belajar', 'NLP.']
```

Berikut adalah penjelasan dari regex yang digunakan:

- `\s+` : Mencocokkan satu atau lebih spasi.

Latihan I Preprocessing Data Teks

Pada latihan ini akan menggunakan dataset dari sebuah project yang bernama GUTENBERG EBOOK. Dataset ini berisi teks dari buku-buku yang sudah di public domain.

Pada latihan ini akan menggunakan dataset berikut:

<https://www.gutenberg.org/files/2638/2638-0.txt>

Jalankan kode berikut untuk mendownload dataset tersebut:

```
1 import re
2 import requests
3
4 def get_gutenberg_book():
5     url = 'https://www.gutenberg.org/files/2638/2638-0.txt'
6
7     raw = requests.get(url).text
8
9     start = re.search(r"\*\*\* START OF THE PROJECT GUTENBERG EBOOK
10                     THE IDIOT \*\*\*", raw).end()
11     stop = re.search(r"II", raw).start()
12
13     text = raw[start:stop]
14
15     return text
16
17 book = get_gutenberg_book()
18 print(book)
```

Selanjutnya adalah mengambil semua kata yang ada pada teks tersebut dan merubahnya menjadi huruf kecil.

```
1 clean_book = re.sub('[^A-Za-z0-9.]+' , ' ', sentence).lower()
2 print(clean_book)
```

Selanjutnya adalah mencoba untuk menghitung jumlah kata 'the' yang ada pada teks tersebut.

```
1 count_the = len(re.findall(r'the', clean_book))
2 print(count_the)
```

Materi 2 - Handling String pada Python

Teknik Preprocessing Data Teks yang kedua adalah Handling String menggunakan Python. Python memiliki berbagai fungsi bawaan yang dapat digunakan untuk melakukan manipulasi string.

Menggabungkan String

```
1 s1 = "Hello"
2 s2 = "World"
3 s3 = s1 + s2
4
5 # Output
6 HelloWorld
```

Mencari sebuah substring dalam sebuah string

```
1 s = "I am learning NLP"
2 print(s.find("learn"))
3
4 # Output
5 5
```

Mengganti sebuah substring dalam sebuah string

```
1 s = "I am learning NLP"
2 print(s.replace("NLP", "Natural Language Processing"))
3
4 # Output
5 I am learning Natural Language Processing
```

Slice String

```
1 String_v1 = "I am exploring NLP"
2 print(String_v1[5:14]) # Memotong string dari index 5 - 14
3 # Output
4 exploring
5
6 print(String_v1[:5]) # Memotong string dari index 0 - 5
7 # Output
8 I am
9
10 print(String_v1[5:]) # Memotong string dari index 5 sampai akhir
11 # Output
12 exploring NLP
13
14 print(String_v1[-3:]) # Memotong string dari index 3 dari akhir
15 # Output
16 NLP
17
18 print(String_v1[:-3]) # Memotong string dari index 0-3 dari akhir
19 # Output
20 I am exploring
```

Untuk lebih lengkapnya, silahkan kunjungi
<https://www.programiz.com/python-programming/methods/string>.

Materi 3 - Ekplorasi dan Preprocessing Data Teks

Pada materi ini akan membahas tentang cara melakukan ekplorasi dan preprocessing data teks menggunakan NLTK dan Pandas. Sebelum mengikuti materi ini, pastikan sudah menginstall NLTK dan Pandas.

Menginstall NLTK dan Pandas

```
pip install nltk
pip install pandas
```

Membaca Data Teks dan Menjadikan DataFrame

```
1 import pandas as pd
2
3 text=[
4     'This is introduction to NLP',
5     'It is likely to be useful, to people ',
6     'Machine learning is the new electrcity',
7     'R is good langauge',
8     'I like this book',
9     'I want more books like this']
10
11 df = pd.DataFrame({'tweet': text})
12 print(df)
```

DataFrame adalah struktur data dua dimensi yang terdiri dari baris dan kolom yang dapat digunakan untuk menyimpan data dalam bentuk tabel. Perintah **pd.DataFrame()** digunakan untuk membuat DataFrame dari data yang diberikan. Silahkan jalankan kode diatas dan perhatikan outputnya.

Menhapus semua tanda baca yang ada pada teks

Dalam preprocessing data teks, seringkali kita perlu menghapus tanda baca yang ada pada teks. Penghapusan ini agar tanda baca tidak mempengaruhi analisis yang akan dilakukan.

```
1 # Contoh
2 s = "I. like. This book!"
3 s1 = re.sub(r'[\w\s]', '', s)
4 print(s1) # 'I like This book'
5
6 # atau bisa dengan menggunakan perulangan
7 import string
8
9 s = "I. like. This book!"
10
11 for c in string.punctuation:
12     s = s.replace(c, "")
13
14 print(s) # 'I like This book'
```

Selanjutnya ubahlah hilangkan semua tanda baca pada teks yang ada pada DataFrame yang sudah dibuat sebelumnya.

```
1 df['tweet'] = df['tweet'].str.replace(r'[^\w\s]', '')
2 print(df['tweet'])
```

Menghapus Stopwords

Stopwords adalah kata-kata yang sering muncul dalam teks dan tidak memiliki makna yang signifikan seperti 'the', 'is', 'and' dan lain-lain atau dalam bahasa Indonesia seperti 'adalah', 'aku', 'apa' dan lain-lain. Sebelum itu silahkan install library NLTK dengan menjalankan perintah berikut untuk mengunduh stopwords yang ada pada NLTK.

```
1 import nltk
2 nltk.download('stopwords')

1 # Sebelum menghapus stopwords, ubahlah teks menjadi huruf kecil
2
3 df['tweet'] = df['tweet'].apply(lambda x: " ".join(x.lower() for x
4         in x.split()))
5
6 from nltk.corpus import stopwords
7
8 stopword = stopwords.words('english')
9 df['tweet'] = df['tweet'].apply(lambda x: " ".join(x for x in x.
10         split() if x not in stop))
11 print(df['tweet'])
```

Standardisasi teks

Standardisasi teks adalah proses mengubah teks menjadi format yang seragam. Misalnya merubah kata-kata yang disingkat menjadi kata lengkap atau kata yang benar.

```
1 # Buatlah sebuah dictionary yang berisi kata-kata yang akan diubah
2 lookup_dict = {
3     'nlp': 'natural language processing',
4     'u': 'you',
5     'r': 'are',
6     'd': 'the',
7     'ur': 'your'}
8
9 def text_std(input_text):
10     words = input_text.split()
11     new_words = []
12     for word in words:
13         word = re.sub(r'[^\w\s]', '', word)
14         if word.lower() in lookup_dict:
15             word = lookup_dict[word.lower()]
16             new_words.append(word)
17         new_text = " ".join(new_words)
18     return new_text
19
```

```
20 print(text_std("I like nlp"))
21
22 # Output
23 natural language processing
```

Spelling Correction (Perbaikan Ejaan)

Spelling correction adalah proses untuk memperbaiki kata-kata yang salah pada pengejaannya. Salah satu library yang dapat digunakan untuk spelling correction adalah **TextBlob** ataupun **autocorrect**. Untuk menggunakan TextBlob dan autocorrect, silahkan install library tersebut dengan menjalankan perintah berikut:

```
pip install textblob
pip install autocorrect
```

```
1 # Menggunakan TextBlob
2 from textblob import TextBlob
3
4 text = "Machine learning is the new electrcity"
5 new_text = TextBlob(text).correct()
6
7 print(new_text) # Output: Machine learning is the new electricity
8
9 # Menggunakan autocorrect
10 from autocorrect import Speller
11 spell = Speller(lang='en')
12
13 text = "Machine learning is the new electrcity"
14 new_text = spell(text)
15
16 print(new_text) # Output: Machine learning is the new electricity
```

Pada contoh diatas, kita menggunakan TextBlob dan autocorrect untuk melakukan spelling correction pada teks yang diberikan. TextBlob akan mengubah kata 'electrcity' menjadi 'electricity' karena kata tersebut salah dalam pengejaannya. Selanjutnya lakukan spelling correction pada teks yang ada pada DataFrame yang sudah dibuat sebelumnya.

```
1 df['tweet'].apply(lambda x: str(TextBlob(x).correct()))
```

Tokenisasi

Tokenisasi adalah proses memecah teks menjadi bagian-bagian yang lebih kecil yang disebut dengan token. Tokenisasi dapat dilakukan dengan menggunakan TextBlob, NLTK atau menggunakan fungsi split() pada Python.

```
1 # Menggunakan TextBlob
2 from textblob import TextBlob
3
4 text = "Machine learning is the new electrcity"
5 new_text = TextBlob(text).words
6
7 print(new_text) # Output: ['Machine', 'learning', 'is', 'the', 'new',
                        ', 'electricity']
```

```
1 # Menggunakan NLTK
2 from nltk.tokenize import word_tokenize
3
4 text = "Machine learning is the new electrcity"
5 new_text = word_tokenize(text)
6
7 print(new_text) # Output: ['Machine', 'learning', 'is', 'the', 'new',
                        ', 'electricity']

1 # Menggunakan split()
2 text = "Machine learning is the new electrcity"
3
4 new_text = text.split()
5 print(new_text) # Output: ['Machine', 'learning', 'is', 'the', 'new',
                        ', 'electricity']
```

Stemming

Stemming adalah proses menghilangkan imbuhan pada kata untuk mendapatkan kata dasar. NLTK memiliki library yang dapat digunakan untuk melakukan stemming yaitu PorterStemmer.

```
1 from nltk.stem import PorterStemmer
2
3 stemmer = PorterStemmer()
4
5 text = "There are many fishes in pound"
6 new_text = " ".join([stemmer.stem(word) for word in text.split()])
7 print(new_text) # Output: there are mani fish in pound
```

Untuk melakukan stemming pada bahasa Indonesia, kita dapat menggunakan library Sastrawi. Penggunaan PySastrawi nanti akan dibahas pada pertemuan selanjutnya.

Lemmatization

Lemmatization adalah proses mengubah kata-kata yang ada pada teks menjadi kata dasar. Lemmatization lebih kompleks dibandingkan dengan stemming karena lemmatization memperhitungkan konteks kata dalam kalimat. Contoh lemmatization adalah kata 'leaves' akan diubah menjadi 'leaf' karena kata 'leaf' adalah kata dasar dari kata 'leaves'.

```
1 import nltk
2 nltk.download('wordnet')

1 # Lemmatization menggunakan TextBlob
2 from textblob import Word
3
4 text = "leaves and leaf"
5 new_text = " ".join([Word(word).lemmatize() for word in text.split()])
6 print(new_text) # Output: leaf and leaf
```



```
1 # Lemmatization menggunakan NLTK
2 from nltk.stem import WordNetLemmatizer
3
4 lemmatizer = WordNetLemmatizer()
5
6 text = "leaves and leaf"
7 new_text = " ".join([lemmatizer.lemmatize(word) for word in text.
8                       split()])
9 print(new_text) # Output: leaf and leaf
```

Tugas

1. Kerjakan seluruh materi yang sudah dijelaskan pada modul diatas.