

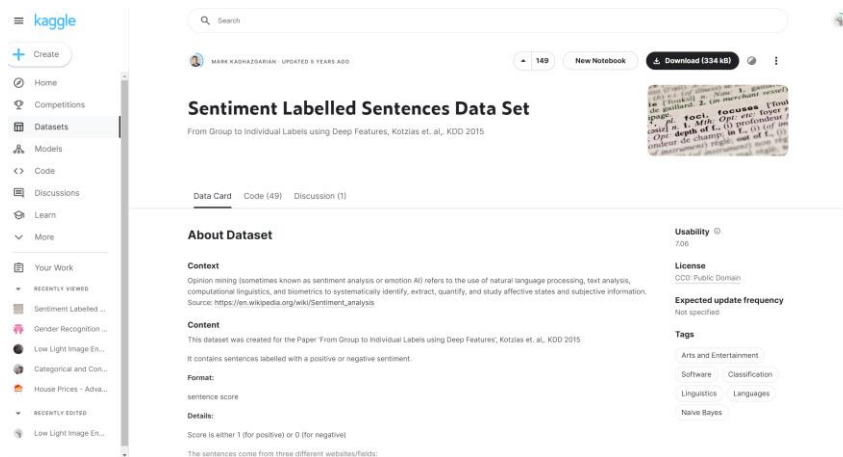
Soal

Untuk data yang ada pada link: <https://www.kaggle.com/marklv/sentiment-labelled-sentences-data-set>, silahkan diolah dengan menggunakan text processing sehingga menghasilkan data bersih. Dikerjakan dengan jupyter notebook pada komputer ataupun google colaboratory, dan upload hasil pengerjaan dalam bentuk ipynb, rar, ataupun pdf (yang penting bisa terbaca oleh penilai).

Note : Untuk mengubah file ipynb ke pdf dapat lebih mudah dilakukan dengan menggunakan google colaboratory.

Link dataset

<https://www.kaggle.com/datasets/marklv/sentiment-labelled-sentences-data-set>



Source Code

```
1 import pandas as pd
2 import re
3 import string
4 import nltk
5 from nltk.corpus import stopwords
6 from nltk.tokenize import word_tokenize
7
8 # File path of the dataset CSV file
9 dataset_file = "C:/Users/LEGION 5 PRO/OneDrive/Documents/Semester 4/Temu Kembali Informasi/dataset 14/sentiment labelled sentences/datasets.csv"
10
11 # Define the regex pattern to remove punctuation and replace with whitespace
12 regex = re.compile('[%s]' % re.escape(string.punctuation.replace('.', ' ')))
13
14 # Define stopwords
15 nltk.download('stopwords')
16 stop_words = set(stopwords.words('english'))
17
18 def preprocess_text(text):
19     # Replace "1", "0", "1,0,0,0", and "1,1,0,0" with whitespace
20     text = re.sub(r'\b[01]\b', ' ', text)
21     text = re.sub(r',', ' ', text)
22     text = re.sub(r',', ' ', text)
23
24     # Tokenize the text
25     tokens = word_tokenize(text)
26
27     # Remove stopwords and punctuation
28     tokens = [token for token in tokens if token.lower() not in stop_words and token not in string.punctuation]
29
30     # Join the tokens back into a string
31     preprocessed_text = ' '.join(tokens)
32
33     return preprocessed_text
34
35 # Read the dataset from the CSV file
36 dataset_files = pd.read_csv(dataset_file).fillna('')
37
38 # Apply preprocessing to each column in the dataframe
39 for column in dataset_files.columns:
40     dataset_files[column] = dataset_files[column].apply(preprocess_text)
41
42 # Save the preprocessed dataset as a new CSV file
43 output_file = "C:/Users/LEGION 5 PRO/OneDrive/Documents/Semester 4/Temu Kembali Informasi/dataset 14/sentiment labelled sentences/preprocessed_dataset.csv"
44 dataset_files.to_csv(output_file, index=False)
45
46 print("Datasets preprocessed and saved to:", output_file)
```

Penjelasan

Kode ini adalah contoh implementasi dari proses pra-pemrosesan teks pada dataset menggunakan library Pandas dan NLTK (Natural Language Toolkit).

Pada awalnya, library yang diperlukan diimpor, seperti Pandas untuk membaca dan menyimpan dataset dalam format CSV, re untuk melakukan operasi pemrosesan teks menggunakan ekspresi reguler, string untuk mendefinisikan kumpulan tanda baca yang akan dihapus dari teks, nltk (Natural Language Toolkit) untuk melakukan pemrosesan teks seperti tokenisasi, dan nltk.corpus untuk mengakses kamus stopwords.

Kemudian, path file dataset ditentukan untuk mengakses file CSV yang berisi dataset. Selanjutnya, pola regex didefinisikan untuk menghapus tanda baca dari teks dan menggantikannya dengan spasi, kecuali tanda baca '.' dan ','.

Pada tahap selanjutnya, kamus stopwords dalam bahasa Inggris diunduh menggunakan `nltk.download('stopwords')`. Stopwords adalah kata-kata yang umumnya tidak memberikan informasi penting dalam pemrosesan teks.

Setelah itu, sebuah fungsi bernama `preprocess_text` dibuat untuk melakukan proses pra-pemrosesan teks. Fungsi ini melakukan beberapa langkah, seperti menggantikan angka "1" dan "0" dengan spasi, menggantikan ",0,," dengan spasi, dan menggantikan ",1,," dengan spasi. Kemudian, teks di-tokenisasi menggunakan `word_tokenize`, yaitu membaginya menjadi token-token kata. Stopwords dan tanda baca dihapus dari token-token tersebut. Akhirnya, token-token digabungkan kembali menjadi string teks.

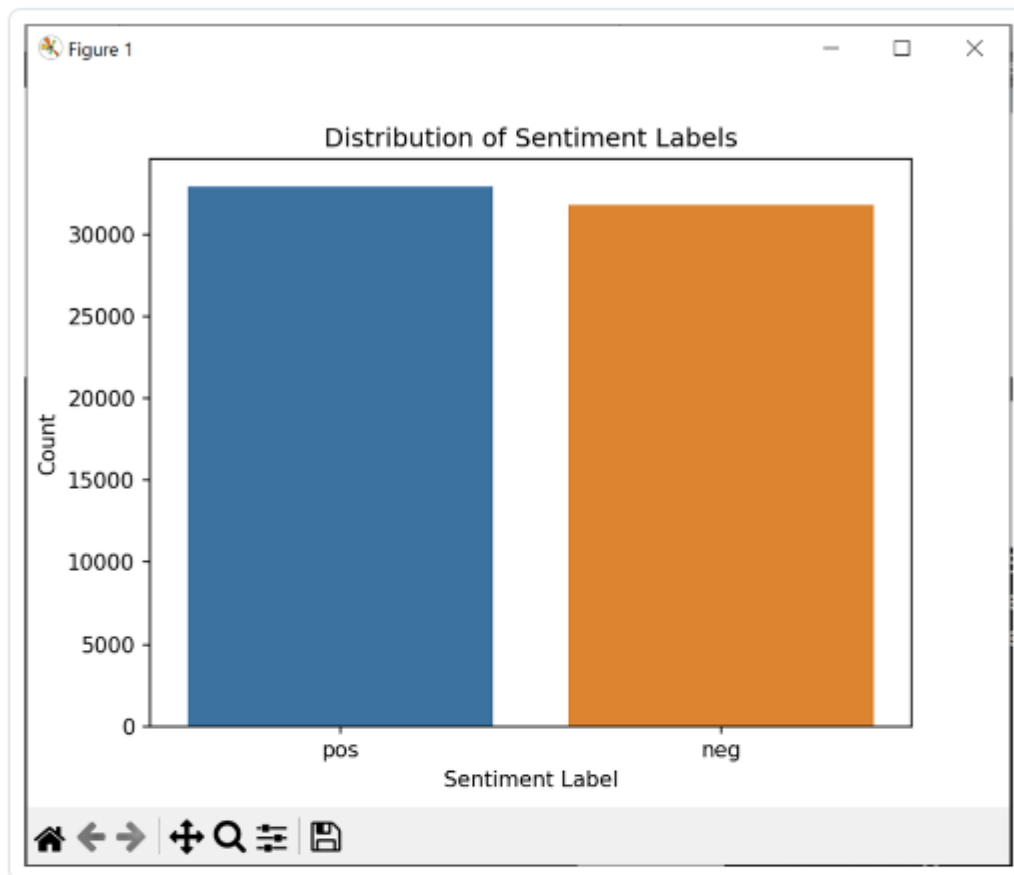
Dataset dibaca dari file CSV menggunakan Pandas, dan nilai-nilai kosong dalam dataset diisi dengan string kosong. Selanjutnya, proses pra-pemrosesan dilakukan pada setiap kolom dalam dataframe. Looping digunakan untuk mengiterasi melalui setiap kolom dataset, dan fungsi `preprocess_text` dipanggil untuk setiap nilai dalam kolom dataset. Hasilnya disimpan dalam kolom yang sama.

Dataset yang telah diproses kemudian disimpan sebagai file CSV baru. Path file output ditentukan, dan dataframe yang telah diproses disimpan ke file CSV baru dengan menggunakan metode `to_csv` dari Pandas. Indeks tidak disertakan dalam file output.

Terakhir, sebuah pesan konfirmasi dicetak untuk menginformasikan bahwa dataset telah diproses dan disimpan dalam file CSV baru. Dengan demikian, kode ini melakukan proses pra-pemrosesan teks pada dataset dengan menghapus stopwords dan tanda baca, dan menyimpan dataset yang telah diproses dalam file CSV baru.

Output





Sentiment analysis in movie_review dataset

Course: Basic Text Processing

Menunggu Approval

Sertifikat

2023/221/19659



CERTIFICATE

OF APPRECIATION

Sertifikat ini diberikan kepada :

HEYDAR EMIR ALVARO

Telah Menyelesaikan Pembelajaran Free Course Mengenai Materi
Basic Text Processing

June 12, 2023

M. Octaviano Pratama S.Kom., M.Kom
President Director Of BISA AI ACADEMY

