

Nama : Virgie Yunita Salsabil

NIM : 21110022

Kelas : S1-SD02A

## ✓ Import Libraries

```
1 # Library untuk operasi input/output
2 import io
3 # Library untuk menghasilkan bilangan acak atau pemilihan elemen acak
4 import random
5 # Library untuk manipulasi string
6 import string
7 # Library untuk memberikan peringatan
8 import warnings
9 # Library untuk operasi array numerik dan manipulasi data
10 import numpy as np
11 # Library untuk pemrosesan teks menggunakan TF-IDF
12 from sklearn.feature_extraction.text import TfidfVectorizer
13 # Library untuk mengukur kemiripan kosinus antara dua set dokumen
14 from sklearn.metrics.pairwise import cosine_similarity
15 # Menyaring peringatan agar tidak ditampilkan
16 warnings.filterwarnings('ignore')
```

```
1 # Instalasi library Natural Language Toolkit (nltk) menggunakan pip
2 !pip install nltk

Requirement already satisfied: nltk in /usr/local/lib/python3.10/dist-packages (3.8.1)
Requirement already satisfied: click in /usr/local/lib/python3.10/dist-packages (from nltk) (8.1.7)
Requirement already satisfied: joblib in /usr/local/lib/python3.10/dist-packages (from nltk) (1.3.2)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.10/dist-packages (from nltk) (2023.6.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk) (4.66.1)
```

```
1 # Import library NLTK
2 import nltk
3 # Import library untuk lemmatization
4 from nltk.stem import WordNetLemmatizer
5 # Pengunduhan data populer dari NLTK
6 nltk.download('popular', quiet=True)
7 #nltk.download('punkt') # first-time use only
8 #nltk.download('wordnet') # first-time use only
```

True

## ✓ Reading in the corpus

Program ini menggunakan halaman Wikipedia sebagai corpus dari chatpotdengan cara meng-Copy konten halaman dan menyimpannya dalam file 'chatbot.txt'

```
1 # Membuka file 'chatbot.txt' untuk dibaca ('r')
2 # Parameter 'errors='ignore'' digunakan untuk mengabaikan karakter yang tidak dapat di-decode
3 f = open('chatbot.txt', 'r', errors='ignore')
4
5 # Membaca seluruh isi file dan menyimpannya dalam variabel 'raw'
6 # raw kini berisi semua data dari corpus per baris (raw)
7 raw = f.read()
8
9 # Mengonversi semua teks dalam 'raw' menjadi huruf kecil (lowercase)
10 raw = raw.lower()
```

## ✓ Tokenisasi

```
1 # Tokenisasi adalah memilah-milah dokumen ke kalimat-kalimat,
2 # kemudian memilah setiap kalimat menjadi sekumpulan kata kata
3 sent_tokens = nltk.sent_tokenize(raw) # converts dokumen corpus ke kalimat-kalimat
4 word_tokens = nltk.word_tokenize(raw) # converts dokumen corpus ke kata-kata
```

## ✓ Preprocessing

```

1 # Membuat objek lemmatizer menggunakan WordNetLemmatizer dari NLTK
2 lemmer = nltk.stem.WordNetLemmatizer()
3
4 # WordNet adalah kamus berorientasi semantik Bahasa Inggris yang disertakan dalam NLTK.
5 # Fungsi LemTokens(tokens) untuk melakukan lemmatization pada setiap token dalam list 'tokens'
6 def LemTokens(tokens):
7     return [lemmer.lemmatize(token) for token in tokens]
8
9 # Membuat kamus yang akan digunakan untuk menghapus tanda baca dari teks
10 remove_punct_dict = dict((ord(punct), None) for punct in string.punctuation)
11
12 # Fungsi LemNormalize(text) untuk melakukan normalisasi teks dengan tokenisasi, lemmatization, dan konversi huruf kecil
13 def LemNormalize(text):
14     return LemTokens(nltk.word_tokenize(text.lower().translate(remove_punct_dict)))

```

## ✓ Keyword matching

```

1 # Daftar kata-kata pembuka yang mungkin diterima sebagai salam
2 GREETING_INPUTS = ("hello", "hi", "greetings", "sup", "what's up", "hey", "hai")
3
4 # Respon-respon yang mungkin diberikan sebagai tanggapan terhadap salam
5 GREETING_RESPONSES = ["hi", "hey", "*nods*", "hi there", "hello", "I am glad! You are talking to me"]
6
7 # Fungsi greeting(sentence) untuk menentukan respon terhadap kalimat yang mengandung kata-kata pembuka
8 def greeting(sentence):
9     # Melakukan iterasi pada setiap kata dalam kalimat
10    for word in sentence.split():
11        # Memeriksa apakah kata dalam bentuk huruf kecil terdapat dalam daftar kata-kata pembuka
12        if word.lower() in GREETING_INPUTS:
13            # Jika ditemukan kata pembuka, maka pilih secara acak salah satu respon dari daftar respon
14            return random.choice(GREETING_RESPONSES)

```

## ✓ Generating Response

After the initial preprocessing phase, we need to transform text into a meaningful vector (or array) of numbers. The bag-of-words is a representation of text that describes the occurrence of words within a document. It involves two things:

- A vocabulary of known words.
- A measure of the presence of known words.

Why is it called a “bag” of words? That is because any information about the order or structure of words in the document is discarded and the model is only concerned with whether the known words occur in the document, not where they occur in the document. The intuition behind the Bag of Words is that documents are similar if they have similar content. Also, we can learn something about the meaning of the document from its content alone. For example, if our dictionary contains the words {Learning, is, the, not, great}, and we want to vectorize the text “Learning is great”, we would have the following vector: (1, 1, 0, 0, 1).

## ✓ TF-IDF Approach

A problem with the Bag of Words approach is that highly frequent words start to dominate in the document (e.g. larger score), but may not contain as much “informational content”. Also, it will give more weight to longer documents than shorter documents. One approach is to rescale the frequency of words by how often they appear in all documents so that the scores for frequent words like “the” that are also frequent across all documents are penalized. This approach to scoring is called Term Frequency-Inverse Document Frequency, or TF-IDF for short, where:

Term Frequency: is a scoring of the frequency of the word in the current document.

- $TF = (\text{Number of times term } t \text{ appears in a document}) / (\text{Number of terms in the document})$

Inverse Document Frequency: is a scoring of how rare the word is across documents.

- $IDF = 1 + \log(N/n)$ , where,  $N$  is the number of documents and  $n$  is the number of documents a term  $t$  has appeared in.

## ✓ Cosine Similarity

Tf-idf weight is a weight often used in information retrieval and text mining. This weight is a statistical measure used to evaluate how important a word is to a document in a collection or corpus

- Cosine Similarity  $(d1, d2) = \text{Dot product}(d1, d2) / \|d1\| * \|d2\|$

where d1,d2 are two non zero vectors. To generate a response from our bot for input questions, the concept of document similarity will be used. We define a function response which searches the user's utterance for one or more known keywords and returns one of several possible responses. If it doesn't find the input matching any of the keywords, it returns a response: "I am sorry! I don't understand you"

```
1 # Fungsi untuk menetapkan respon jawaban
2 def response(user_response):
3     robo_response = '' # Pada tahap awal, respon mesin diisi karakter kosong
4     sent_tokens.append(user_response) # Pertanyaan user ditokenisasi dan ditambahkan ke corpus
5
6     # Posisi paling bawah (yaitu posisi -1)
7     TfidfVec = TfidfVectorizer(tokenizer=LemNormalize, stop_words='english')
8     tfidf = TfidfVec.fit_transform(sent_tokens) # Token dari pertanyaan user di-vektorisasi
9     vals = cosine_similarity(tfidf[-1], tfidf) # Hitung similarity setiap token corpus dengan token pertanyaan
10
11     # Sortir jarak similarity setiap token corpus dengan token pertanyaan
12     # Mengambil indeks token yang memiliki similarity tertinggi (indeks paling rendah) dengan pertanyaan pengguna
13     idx = vals.argsort()[0][-2]
14
15     # Meratakan dan mengurutkan nilai similarity untuk mendapatkan similarity tertinggi kedua
16     flat = vals.flatten()
17     flat.sort()
18
19     # Menyimpan nilai similarity tertinggi kedua
20     req_tfidf = flat[-2]
21
22     # Jika pertanyaan dan semua token corpus jaraknya tinggi, maka pertanyaan tidak memiliki jawaban
23     if req_tfidf == 0:
24         robo_response = robo_response + "Please reply, I don't understand your questions"
25         return robo_response
26     else:
27         # Jika jaraknya terendah, maka token tersebut dipakai sebagai jawaban
28         robo_response = robo_response + sent_tokens[idx]
29         return robo_response
```

Finally, we will feed the lines that we want our bot to say while starting and ending a conversation depending upon user's input.

```
1 # Inisialisasi variabel 'flag' untuk menjalankan loop
2 flag=True
3 print("Mesin: My name is Robo. I will answer your queries about Chatbots. If you want to exit, type Bye!")
4
5 # Memulai loop utama, akan terus berjalan selama 'flag' bernilai True
6 while(flag==True):
7     # Meminta user untuk memasukkan pertanyaan
8     user_response = input("Masukkan pertanyaan :")
9     # Mengubah input pengguna menjadi huruf kecil untuk konsistensi
10    user_response=user_response.lower()
11
12    # Memeriksa jika user tidak keluar
13    if(user_response!='bye'):
14        # Jika pengguna mengucapkan thanks/thankyou, menghentikan loop & memberikan balasan
15        if(user_response=='thanks' or user_response=='thank you') :
16            flag=False #tandai proses berhenti
17            print("Mesin: You are welcome..") #balasan thank you
18        else:
19            if(greeting(user_response)!=None): #jika response adalah kalimat greeting
20                print("Mesin: "+greeting(user_response)) #tampilkan kalimat greeting
21            else: #jika bukan kalimat greeting
22                print("Mesin: ",end="")
23                print(response(user_response)) #memproses user answer disini
24                sent_tokens.remove(user_response) #user answer dihapus setelah dimunculkan
25                print("-----")
26
27    # Jika pengguna ingin keluar, mengubah nilai flag menjadi False untuk menghentikan loop
28    else:
29        flag=False
30        # Menampilkan pesan perpisahan dan pemisah untuk setiap iterasi
31        print("Mesin: Bye! take care..")
32        print("=====")
```

```
Mesin: My name is Robo. I will answer your queries about Chatbots. If you want to exit, type Bye!
Masukkan pertanyaan :hi
Mesin: hi there
Masukkan pertanyaan :What are the categories of chatbot usage?
Mesin: analytics
```

the usage of the chatbot can be monitored in order to spot potential flaws or problems.

Masukkan pertanyaan :Who coined the term "ChatterBot" and when?

Mesin: the term "chatterbot" was originally coined by michael mauldin (creator of the first verbot, julia) in 1994 to describe these

Masukkan pertanyaan :What is the difference between older chatbots like ELIZA and PARRY compared to newer generations?

Mesin: while eliza and parry were used exclusively to simulate typed conversation, many chatbots now include functional features suc

Masukkan pertanyaan :How is the chatbot development process structured, and what are its key phases?

Mesin: chatbot creation

the process of creating a chatbot follows a pattern similar to the development of a web page or a mobile app.

Masukkan pertanyaan :What is the impact of artificial intelligence (AI) usage on jobs, according to Forrester?

Mesin: according to forrester (2015), ai will replace 16 percent of american jobs by the end of the decade.chatbots have been used i

Masukkan pertanyaan :thank you

Mesin: You are welcome..

## KESIMPULAN

- Berdasarkan output di atas, fungsi **cosine similarity** dalam program berjalan dengan baik. Ketika pengguna memasukkan pertanyaan, program menghitung kemiripan cosinus antara vektor TF-IDF dari pertanyaan tersebut dengan vektor TF-IDF dari seluruh korpus (pertanyaan sebelumnya). Hasil similarity digunakan untuk menentukan jawaban terbaik dengan memilih pertanyaan sebelumnya yang paling mirip.
- Contohnya, ketika pengguna memasukkan pertanyaan "What are the categories of chatbot usage?" program memberikan jawaban yang relevan dengan mencocokkan kata kunci dan struktur pertanyaan dengan pertanyaan yang ada dalam korpus. Ini menunjukkan bahwa cosine similarity berfungsi dengan baik dalam memahami konteks pertanyaan dan memberikan respons yang sesuai.