

webSummerizer

Hong Huy Hoang

November 2024

1 Introduction

My second project using LLM, using Google's Gemini generative AI

2 Requirements

- language Python 3.11.10
- beautifulsoup4 4.12.2
- ipython 8.15.0
- requests 2.31.0
- google-generative-ai 0.3.0

3 Setup

I recommend using miniconda, although every other python virtual env managers will work fine

The command below will install a conda environment. You can change that of the environment in environment.yml before running it

```
1 conda env create -f environment.yml
2 conda activate websummerizer
```

4 How this works?

4.1 Load Environment and API from Google AI Studio

```
1 #load env
2 import requests
3 import google.generativeai as genai
```

```

4 import os
5 from bs4 import BeautifulSoup
6 from ipykernel.kernelapp import IPKernelApp
7 from IPython.display import Markdown, display

```

```

1 #load API
2 os.environ["GOOGLE_API_KEY"] = "your_api_here"
3 genai.configure(api_key="your_api_here")
4 generation_config = {"temperature" : 0.9, "top_p" : 1, "top_k" : 1}
5 model = genai.GenerativeModel("gemini-pro", generation_config=generation_config)

```

4.2 Define a class represent a webpage

```

1 #define a class represent a webpage
2 #beautifulSoup is used for web scanning and extract information from html
3
4 class Website:
5     url : str
6     title : str
7     text : str
8
9     def __init__(self, url):
10         self.url = url
11         response = requests.get(url)
12         soup = BeautifulSoup(response.text, 'html.parser')
13         self.title = soup.title.string if soup.title else "No title found"
14
15         #remove unnecessary things
16         for irrelevant in soup.body(["script", "style", "img", "input"]):
17             irrelevant.decompose()
18
19         #strip text
20         self.text = soup.body.get_text(separator = "\n", strip = True)
21

```

4.3 Define system prompt and user prompt for AI to know more about context

System prompt

```

1 system_prompt = "You are an assistant that analyzes the contents of a website \
2     and provides a short summary, ignoring text that might be navigation related. \

```

```
3     Respond in markdown."
```

```
4
```

User prompt

```
1     def user_prompt_for(website):
2         user_prompt = f"You are looking at a website titled {website.title}"
3         user_prompt += "The contents of this website is as follows; \
4             please provide a short summary of this website in markdown.\
5             If it includes news or announcements, then summerize these too \n\n"
6         user_prompt += website.text
7     return user_prompt
```

4.4 Time to bring all together, very simple!

```
1     def summarize(url):
2         website = Website(url)
3         full_prompt = f"System: {system_prompt} User: {user_prompt_for(website)}"
4         response = model.generate_content(full_prompt)
5     return response.text
```

4.5 Testing Gemini AI model

Example url: <https://daa.uit.edu.vn/>

```
1     summarize('https://daa.uit.edu.vn/')
2     display_summary('https://daa.uit.edu.vn/')
```

Output would be something like this



Figure 1: Example of output