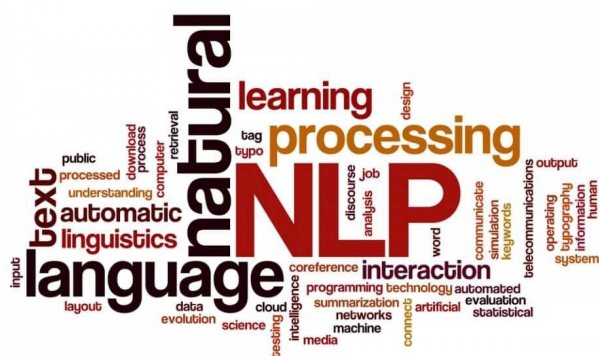


بسم الله الرحمن الرحيم



دانشکده مهندسی برق و کامپیوتر

پردازش زبان‌های طبیعی - تمرین ششم

سید مهدی رضوی

استاد : آقای دکتر فیلی

تیر ماه ۱۴۰۳

فهرست مطالب

۲	۱ بخش اول
۳	۲ بخش دوم
۵	۳ بخش سوم
۶	۴ بخش چهارم

فهرست تصاویر

۳	۱ بارگذاری تمام لینک‌ها در صفحه رسمی کتاب
۴	۲ لینک فصل‌های مختلف کتاب ژورافسکی
۴	۳ عملیات قطعه‌بندی کردن متون
۵	۴ اسناد بازگردانده‌شده توسط بازیاب در جواب یک پرس‌وجو مرتبط با پردازش زبان
۶	۵ اسناد بازگردانده‌شده توسط بازیاب در جواب یک پرس‌وجو مرتبط با پردازش زبان
۶	۶ اسناد بازگردانده‌شده توسط بازیاب در جواب یک پرس‌وجو مرتبط با علوم کامپیوتر
۷	۷ سناد بازگردانده‌شده توسط بازیاب در جواب یک پرس‌وجو غیرمرتبط
۷	۸ ChatPromptTemplate
۷	۹ ChatPromptTemplate
۸	۱۰ llm invoke

۱ بخش اول

RecursiveCharacterTextSplitter یک ابزار کاربردی در LangChain است که برای تقسیم اسناد طولانی به قطعات کوچکتر و قابل مدیریت طراحی شده است. این به ویژه در هنگام برخورد با متون بزرگی که نیاز به پردازش، نمایه سازی یا جستجو دارند، مفید است. در اینجا توضیح مفصلی در مورد chunksize و chunkoverlap آمده است:

تعریف: chunksize حداکثر تعداد کاراکترهایی را که هر تکه شامل می شود را مشخص می کند. هدف: این پارامتر تضمین می کند که هر تکه از متن خیلی بزرگ نیست، که به دلایل مختلفی می تواند مفید باشد: عملکرد: تکه‌های کوچک‌تر پردازش آسان‌تر بوده و به حافظه کمتری نیاز دارند. ارتباط: به حفظ تمرکز و زمینه متن در محدوده‌های قابل مدیریت کمک می‌کند و کار مدل‌های زبان و الگوریتم‌های جستجو را آسان‌تر می‌کند. نمایه سازی: تکه های کوچکتر برای نمایه سازی و بازیابی در پایگاه های داده برداری مانند FAISS کارآمدتر هستند.

تعریف: chunkoverlap تعداد کاراکترهایی را که باید بین تکه های متوالی همپوشانی داشته باشند را مشخص می کند. هدف: این پارامتر تضمین می‌کند که بین تکه‌ها همپوشانی وجود دارد که به دلایل مختلفی می‌تواند مفید باشد: Preservation: Context به حفظ زمینه بین تکه ها کمک می کند. این مهم است زیرا برخی از اطلاعات مرتبط ممکن است در مرز یک تکه قطع شوند. تداوم: تضمین می‌کند که یک انتقال آرام بین تکه‌ها وجود دارد، و خطر از دست دادن زمینه معنادار را که ممکن است بین تکه‌ها تقسیم شود، کاهش می‌دهد.

✓ Loading Data

```
import requests
from bs4 import BeautifulSoup

url = "https://stanford.edu/~jurafsky/slp3/"
response = requests.get(url)
soup = BeautifulSoup(response.text, 'html.parser')

pdf_links = []
for link in soup.find_all('a', href=True):
    href = link['href']
    if href.endswith('.pdf'):
        pdf_links.append(url + href if href.startswith('/') else href)

[ ] pdf_links_main = []
for pdf in pdf_links:
    item = 'https://' + append_text + pdf
    pdf_links_main.append(item)
    print(item)
```

Executing (2m 33s) <cell line: 5> > from_documents() > from_texts() > embed_documents() > encode

شکل ۱: بارگذاری تمام لینک‌ها در صفحه رسمی کتاب

۲ بخش دوم

در LangChain، Embedder یک نقش کلیدی، در تبدیل متن به بردارهای عددی ایفا می‌کند. این بردار عددی که معنای متن را به تصویر می‌کشد. این بردارهای عددی برای کارهای مختلف مانند جستجوی شباهت، خوشه‌بندی و بازیابی معنایی ضروری هستند. Embedder سنگ بنای درک معنایی در LangChain است که سیستم را قادر می‌سازد از تطابق واژگانی فراتر رفته و از روابط معنایی عمیق بین متون استفاده کند. این منجر به نتایج دقیق‌تر، معنادارتر و مرتبط‌تر در کاربردهای مختلف می‌شود.

```
urls = [  
    'https://stanford.edu/~jurafsky/slp3/2.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/3.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/4.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/5.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/6.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/7.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/8.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/9.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/10.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/11.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/12.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/13.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/14.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/15.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/16.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/17.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/18.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/19.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/20.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/21.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/22.pdf' ,  
    'https://stanford.edu/~jurafsky/slp3/23.pdf' ,  
]
```

شکل ۲: لینک فصل‌های مختلف کتاب ژورافسکی

✓ Chunking, Vector Store and Retriever

```
text_splitter = RecursiveCharacterTextSplitter(chunk_size = 1024 , chunk_overlap = 64)  
chunks = text_splitter.split_documents(docs)  
  
embedding_function = HuggingFaceEmbeddings(show_progress = True, multi_process = True)  
vector_store = FAISS.from_documents(documents=chunks, embedding = embedding_function)  
  
retriever = vector_store.as_retriever(search_kwargs = {"k": 3})  
  
*** /usr/local/lib/python3.10/dist-packages/langchain_core/_api/deprecation.py:139: LangChainDeprecationWarning: The class `HuggingFaceEmbeddings` is deprecated.  
warn_deprecated(  
/usr/local/lib/python3.10/dist-packages/sentence_transformers/cross_encoder/CrossEncoder.py:11: TqdmExperimentalWarning: Using `tqdm` in a loop can be slow.  
from tqdm.autonotebook import tqdm, trange  
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:89: UserWarning:  
The secret `HF_TOKEN` does not exist in your Colab secrets.  
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret `HF_TOKEN` in your Colab secrets, and restart this notebook.  
You will be able to reuse this secret in all of your notebooks.  
Please note that authentication is recommended but still optional to access public models or datasets.
```

شکل ۳: عملیات قطعه‌بندی کردن متون

- رویکرد بازیاب معنایی بر مبنای Embedding-based matching
- می تواند مفهوم و مفهوم کلمات و جملات را درک کند.
- مترادف ها، نقل قول ها، و تغییرات در جمله بندی را به طور موثر مدیریت می کند.
- برای پرس و جوهای پیچیده که نیاز به درک معنایی دارند، موثرتر است.
- عموماً به دلیل هزینه محاسباتی تولید Embedding ، کندتر از بازیابی های واژگانی است.
- به مدل های از پیش آموزش دیده و گاهی اوقات منابع محاسباتی بزرگ نیاز دارد.
- برای دستیابی به عملکرد مطلوب ممکن است برای موارد استفاده خاص نیاز به تنظیم دقیق داشته باشد.

- رویکرد بازیاب واژگانی
- رویکرد: تطبیق کلمات کلیدی (Algorithms like TF-IDF)
- سریع و کارآمد برای پرس و جوهای ساده.
- زمانی که تطبیق دقیق کلمه کلیدی بسیار مهم است، به خوبی کار می کند.
- نیازی به آموزش یا تعبیه های از پیش محاسبه شده ندارد.
- نقاط ضعف: قادر به درک مفهوم یا مفهوم کلمات نیست.
- نمی تواند مترادف ها یا نقل قول ها را به طور موثر مدیریت کند.
- برای پرس و جوهای پیچیده که در آن معناشناسی اهمیت دارد، کمتر موثر است.

[illegible]

شکل ۴: اسناد بازگردانده شده توسط بازیاب در جواب یک پرس و جو مرتبط با پردازش زبان

شکل ۵: اسناد بازگردانده شده توسط بازیاب در جواب یک پرس و جو مرتبط با پردازش زبان

شکل ۶: اسناد بازگردانده شده توسط بازیاب در جواب یک پرس و جو مرتبط با علوم کامپیوتر

شکل ۷: سناد بازگردانده شده توسط بازتاب در جواب یک پرسش و جو غیر مرتبط

ChatPromptTemplate : شکل ۸

- Section 4

شکل ۹: ChatPromptTemplate

```
from pprint import pprint

# raw_prompt = "Hi, I'm learning langchain to use it with LLMs and make awesome stuff!"
# raw_prompt = "سلام . من علاقه مند به یادگیری تکنولوژی اسپارک و هادوب و کافکا برای یافتن شغل مهندس داده هستم"

raw_prompt2 = 'Hi my dear assistant ! What is NLP and What does it do in our lifes ?'
raw_prompt3 = 'What is NLP Reference book and What\'s the topics of this book ? please tell me in persian . thanks'
llm.invoke(raw_prompt3).content

'Here is the information you requested:\n\n**NLP Reference Book:**\n\nOne of the most popular and widely used NLP reference books is "Natural Language Processing (almost) from Scratch" by Collobert et al. (2011). However, there are many other excellent books on NLP that can serve as a reference, depending on your specific needs and goals.\n\n**Topics Covered in a Typical NLP Reference Book:**\n\nHere are some of the typical topics covered in a comprehensive NLP reference book:\n\n1. **Mathematical Foundations:** Linear Algebra, Calculus, Probability Theory, and Information Theory.\n2. **Language Modeling:** N-gram models, Markov models, and neural network-based language models.\n3. **Text Preprocessing:** Tokenization, Stopword removal, Stemming, Lemmatization, and Named Entity Recognition.\n4. **Word Embeddings:** Word2Vec, GloVe, and FastText.\n5. **Text Classification:** Sentiment Analysis, Spam Detection, and Topic Modeling.\n6. **Sequence Labeling:** Part-of-Speech Tagging, Named Entity Recognition, etc.'
```

```
[ ] prompt_template = "Hi, I'm learning {tool} to use it for classify of {query} in one of the 3 categories. \nNatural Language Processing , Computer Science or None of them."

prompt_chat_template = ChatPromptTemplate.from_template(
    template=prompt_template,
)
```

شکل ۱۰: llm invoke