

به نام خدا



دانشگاه تهران

دانشکده فنی

دانشکده مهندسی برق و کامپیوتر



درس پردازش زبان طبیعی

تمرین ۶

نام و نام خانودگی: علی عطاءاللهی

شماره دانشجویی: ۸۱۰۱۹۹۴۶۱

اسفند ماه ۱۴۰۲

۳ پاسخ سوال اول
۳ دادگان
۳ بخش اول - دریافت و آماده‌سازی دادگان
۵ بخش دوم - تولید و بازنمایی پایگاه داده برداری
۶ بخش سوم - پیاده‌سازی بازیاب معنایی
۸ بخش سوم - پیاده‌سازی بازیاب ترکیبی (امتیازی)
۹ بخش چهارم - پیاده‌سازی Router chain
۱۲ بخش پنجم - پیاده‌سازی search engine chain
۱۴ بخش ششم - پیاده‌سازی relevancy check chain (امتیازی)
۱۷ بخش هفتم - پیاده‌سازی Fallback chain
۲۰ بخش هشتم - پیاده‌سازی Generative with context chain
۲۱ بخش نهم - آماده‌سازی گراف با استفاده از Langgraph



```
1 documents = []
2 for link in links:
3     online_pdf_loader = OnlinePDFLoader(link)
4     documents.extend(online_pdf_loader.load())
```

پ) سپس این مستندات با استفاده از `RecursiveCharacterTextSplitter` به اندازه‌های مناسب تقسیم می‌شوند. این کلاس متن را بر اساس اندازه‌ها و هم‌پوشانی‌های مشخص شده تقسیم می‌کند. برای مثال، با اندازه بلوک ۱۰۲۴ و هم‌پوشانی بلوک ۶۴، متن به بلوک‌های هم‌پوشانده ۱۰۲۴ کاراکتری تقسیم می‌شود که هر بلوک ۶۴ کاراکتر با بلوک قبلی هم‌پوشانی دارد.



```
1 recursive_character_text_splitter = RecursiveCharacterTextSplitter(
2     chunk_size=1024,
3     chunk_overlap=64,
4     length_function=len
5 )
6
7 chunks = recursive_character_text_splitter.split_documents(documents)
```

ت) تقسیم‌بندی متن به بخش‌های کوچکتر به چند دلیل مهم است:

پردازش کارآمد: بلوک‌های کوچکتر آسان‌تر پردازش و ایندکس می‌شوند، به ویژه برای پایگاه‌های داده برداری.

بهبود بازیابی: این کار دقت بازیابی سند را افزایش می‌دهد، زیرا اطلاعات مربوطه در بخش‌های قابل مدیریت‌تر قرار می‌گیرند.

مقیاس‌پذیری: این کار مقیاس‌پذیری را پشتیبانی می‌کند و به سیستم امکان می‌دهد اسناد بزرگ را در بخش‌های کوچکتر و قابل مدیریت‌تر اداره کند.

ث) انتخاب مقادیر مناسب برای اندازه بلوک و هم‌پوشانی حیاتی است:

اندازه بلوک: اگر خیلی بزرگ باشد، بلوک‌ها ممکن است حاوی اطلاعات زیادی باشند که دقت بازیابی را کاهش می‌دهد. اگر خیلی کوچک باشد، ممکن است زمینه از دست برود.

هم‌پوشانی بلوک: اگر خیلی زیاد باشد، ممکن است منجر به تکرار و ناکارآمدی شود. اگر خیلی کم باشد، ممکن است زمینه مهم بین بلوک‌ها از دست برود.

مقادیر پیشنهادی ۱۰۲۴ برای اندازه تکه و ۶۴ برای همپوشانی تکه برای متعادل کردن حفظ بافت و کارایی محاسباتی در وظایف پردازش متن حیاتی هستند. اندازه تکه‌ای از ۱۰۲۴ توکن تضمین می‌کند که مدل دارای زمینه کافی برای درک و تولید متن منسجم است، که برای کارهایی مانند خلاصه‌سازی و پاسخ‌گویی به سؤال ضروری است. همپوشانی ۶۴ توکن، پیوستگی متنی بین تکه‌ها را حفظ می‌کند، و از از دست رفتن اطلاعات مهم بدون ایجاد افزونگی بیش از حد جلوگیری می‌کند. این مقادیر به حفظ کارایی و اثربخشی کمک می‌کنند و تضمین می‌کنند که مدل بدون تجاوز از محدودیت‌های حافظه یا پردازش داده‌های غیرضروری به خوبی عمل می‌کند.

بخش دوم – تولید و بازنمایی پایگاه داده برداری

الف) مستندات تقسیم‌شده با استفاده از HuggingFaceEmbeddings تعبیه شده و در یک پایگاه داده برداری FAISS ذخیره می‌شوند. این مرحله شامل تبدیل متن به بردارهای عددی است که معنای ضمنی را ثبت می‌کنند و جستجوهای شباهتی کارآمد را تسهیل می‌کنند.



```
1 def load_faiss_vector_store(embeddings_db_name):
2     if not os.path.isfile(embeddings_db_name):
3         faiss_vector_store = FAISS.from_documents(chunks, embeddings)
4         faiss_vector_store.save_local(embeddings_db_name)
5     else:
6         faiss_vector_store = FAISS.load_local(embeddings_db_name)
7
8     return faiss_vector_store
```



```
1 store = LocalFileStore('./cache/')
2 general_embedding = HuggingFaceEmbeddings()
3 embeddings = CacheBackedEmbeddings.from_bytes_store(
4     general_embedding,
5     store
6 )
```



```
1 faiss_vector_store = load_faiss_vector_store(  
2     embeddings_db_name="faiss_index"  
3 )  
4  
5 print(faiss_vector_store.index.ntotal)
```

خروجی: ۲۱۲۱

(ب) استفاده از تعبیه‌گر مناسب حیاتی است:

سازگاری زبانی: یک تعبیه‌گر که بر روی زبان یا زمینه متفاوت آموزش دیده باشد ممکن است نتواند ظرافت‌های زبان هدف را ثبت کند و منجر به عملکرد ضعیف بازیابی شود.

بهینه‌سازی عملکرد: مدل‌های تعبیه‌ای که دارای دقت بالا برای زبان و حوزه هدف هستند، نمایه‌سازی و بازیابی بهتری را تضمین می‌کنند.

کشینگ تعبیه‌ها می‌تواند کارایی را بهبود بخشد و امکان استفاده مجدد از تعبیه‌های پیش‌محاسبه شده و صرفه‌جویی در منابع محاسباتی را فراهم کند.

بخش سوم - پیاده‌سازی بازیاب معنایی

الف) پیاده‌سازی شامل استفاده از ابزار FAISS، ادغام شده با چارچوب LangChain، برای تسهیل قابلیت‌های جستجوی معنایی کارآمد و موثر است. به‌ویژه برای مدیریت مجموعه داده‌های بزرگ به دلیل الگوریتم‌های جستجو و ساختارهای شاخص بهینه‌سازی شده‌اش مناسب است.

مفهوم بازی ابری حول ایده شاردینگ می‌چرخد، جایی که مجموعه داده به بخش‌های کوچکتر و قابل مدیریت تقسیم می‌شود. این رویکرد امکان پردازش موازی و زمان‌های بازیابی سریع‌تر را فراهم می‌کند. با پیاده‌سازی FAISS با شاردینگ، می‌توانیم اطمینان حاصل کنیم که فرآیند بازیابی حتی با رشد مجموعه داده‌ها نیز مقیاس‌پذیر و پاسخگو باقی می‌ماند.



```
1 faiss_retriever.invoke(related_query)
4 )
```



```
1 bm25_term_saturation = 4
2 bm25_retriever = BM25Retriever.from_texts(
3     [chunk.page_content for chunk in chunks]
4 )
5
6 bm25_retriever.k = bm25_term_saturation
```

ب) برای ارزیابی عملکرد بازیابی کننده خود، از سه پرسش مختلف استفاده کردیم. با تحلیل مستندات بازیابی شده برای هر پرسش، می‌توانیم توانایی بازیابی کننده در مدیریت انواع مختلف درخواست‌های اطلاعاتی و اثربخشی آن در ارائه نتایج مرتبط را ارزیابی کنیم.



```
1 related_query = "Named Entity Recognizers"
2 unrelated_query = "Cloud Computing"
3 super_unrelated_query = "Champion of Asian cup"
```

```
✓ 0s [14] print(faiss_retriever.invoke(related_query))
      print(faiss_retriever.invoke(unrelated_query))
      print(faiss_retriever.invoke(super_unrelated_query))

[Document(metadata={'source': '/tmp/tmp27_1z6/tmp.pdf'}, page_content=
[Document(metadata={'source': '/tmp/tmp9ck7eis0/tmp.pdf'}, page_content=
[Document(metadata={'source': '/tmp/tmptn84r0wh/tmp.pdf'}, page_content=

✓ 0s print(bm25_retriever.invoke(related_query))
      print(bm25_retriever.invoke(unrelated_query))
      print(bm25_retriever.invoke(super_unrelated_query))

[Document(page_content='8.5.2 Features for CRF Named Entity Recognizers\
[Document(page_content='(cid:80)K\n\n(where c is the correct class)\n\n(
[Document(page_content='9.8\n\nimpression 5.95 3.65 0.98\n\nagreement 0.
```

بخش سوم - پیاده سازی بازیاب ترکیبی (امتیازی)

الف) بازیابی معنایی بر درک معنای پشت کلمات و عبارات در یک پرسش تمرکز دارد. این نوع بازیابی از مدل‌های پیشرفته‌ای مانند تعبیه‌ها (Embeddings) برای ثبت ظرافت‌های متنی و روابط معنایی بین اصطلاحات استفاده می‌کند. این نوع بازیابی به‌ویژه برای مدیریت پرسش‌های پیچیده‌ای که مطابقت دقیق کلمات برای یافتن مستندات مرتبط کافی نیست، موثر است.

بازیابی لغوی، از طرف دیگر، بر تطبیق دقیق کلمات یا عبارات در پرسش با آن‌هایی که در مستندات هستند متکی است. BM25 یک الگوریتم معروف در این دسته است که مستندات را بر اساس فراوانی اصطلاح و فراوانی معکوس اصطلاح امتیازدهی می‌کند. بازیابی لغوی برای پرسش‌های دقیق که در آن اصطلاحات خاص برای بازیابی مستندات مرتبط حیاتی هستند، موثر است.

ترکیب این دو روش به یک سیستم بازیابی جامع‌تر اجازه می‌دهد که بتواند طیف وسیعی از پرسش‌ها را به طور موثر مدیریت کند.

ب) بازیابی‌کننده ترکیبی روش‌های بازیابی معنایی (FAISS) و لغوی (BM25) را با استفاده از کلاس EnsembleRetriever در کتابخانه LangChain ترکیب می‌کند.



```
1 def build_ensemble_retriever(wights):
2     return EnsembleRetriever(
3         retrievers=[bm25_retriever, faiss_retriever], wights=wights
4     )
5
6 def test_ensemble_retriever(wights):
7     ensemble_retriever = build_ensemble_retriever(wights)
8     print(f"for {wights}:")
9     print(ensemble_retriever.invoke(related_query))
10    print(ensemble_retriever.invoke(unrelated_query))
11    print(ensemble_retriever.invoke(super_unrelated_query))
12    print("\n")
```

پ) برای تعیین وزن‌های مناسب برای بازیابی‌کننده ترکیبی، ترکیبات مختلف وزن‌ها برای بازیابی‌کننده‌های BM25 و FAISS را آزمایش کردیم. به‌طور کلی، برای زبان‌هایی مانند انگلیسی که در آن مدل‌های معنایی قدرتمند هستند، وزن بیشتر برای بازیابی معنایی (FAISS) معمولاً منجر به عملکرد بهتر می‌شود. با این حال، تعادل این وزن‌ها برای اطمینان از اینکه هر دو روش به‌طور موثر در فرآیند بازیابی نهایی مشارکت دارند، حیاتی است (در بخش بعدی وزن‌ها مشخص شده‌اند).

ت) عملکرد بازیابی‌کننده ترکیبی را با استفاده از همان سه پرسش ذکر شده در بالا ارزیابی کردیم. نتایج نشان داد که بازیابی‌کننده ترکیبی به‌طور موثری نقاط قوت هر دو روش بازیابی معنایی و لغوی را ترکیب می‌کند و مستندات مرتبط و دقیقی را برای طیف وسیعی از پرسش‌ها ارائه می‌دهد.

```
test_ensemble_retriever([0.5, 0.5])
test_ensemble_retriever([0.1, 0.9])
test_ensemble_retriever([0.9, 0.1])
```

[17]

```
... for [0.5, 0.5]:
[Document(page_content='8.5.2 Features for CRF Named Entity Recognizers\n\nA CRF for NER makes use of very similar features to
[Document(page_content='(cid:80)K\n\n(where c is the correct class)\n\n(7.26)\n\n7.5.2 Computing the Gradient\n\nHow do we com
[Document(page_content='9.8\n\nimpression 5.95 3.65 0.98\n\nagreement 0.3\n\nWord Relatedness The meaning of two words can be

for [0.1, 0.9]:
[Document(page_content='8.5.2 Features for CRF Named Entity Recognizers\n\nA CRF for NER makes use of very similar features to
[Document(page_content='(cid:80)K\n\n(where c is the correct class)\n\n(7.26)\n\n7.5.2 Computing the Gradient\n\nHow do we com
[Document(page_content='9.8\n\nimpression 5.95 3.65 0.98\n\nagreement 0.3\n\nWord Relatedness The meaning of two words can be

for [0.9, 0.1]:
[Document(page_content='8.5.2 Features for CRF Named Entity Recognizers\n\nA CRF for NER makes use of very similar features to
[Document(page_content='(cid:80)K\n\n(where c is the correct class)\n\n(7.26)\n\n7.5.2 Computing the Gradient\n\nHow do we com
[Document(page_content='9.8\n\nimpression 5.95 3.65 0.98\n\nagreement 0.3\n\nWord Relatedness The meaning of two words can be
```

بخش چهارم — پیاده‌سازی ROUTER CHAIN

الف) ابتدا، ما از سایت TogetherAI کلید KeyAPI رایگان خود را با ثبت نام دریافت کردیم. این کلید به ما اعتبار اولیه ارائه می دهد که امکان استفاده از مدل های منبع باز، به خصوص مدل Llama-3-70b ، را بدون محدودیت هزینه ممکن می سازد.

```
1 os.environ["TAVILY_API_KEY"] = "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"
2 os.environ["TOGETHER_API_KEY"] = "9f24f1264668e13d395b5f214fce62f405544086d879b3fa954b5578de15b72d"
```

ب) ما مدل خود را با استفاده از کلاس ChatTogether در LangChain ایجاد کردیم، با استفاده از مدل meta-llama/Llama-3-70b. تنظیم این مدل با دمای 0/0 انجام شد که برای اطمینان از پاسخ های قطعی از مدل بدون متغیریت ورودی استفاده می شود.

```
1 router_llm = ChatTogether(
2     together_api_key=os.environ["TOGETHER_API_KEY"],
3     model="meta-llama/Llama-3-70b-chat-hf",
4     temperature=0.0
5 )
```

پ) زنجیره روتر برای مدیریت وظیفه دسته بندی طراحی شد. این زنجیره از سه قسمت ترکیب شده است:

الگوی پرسش چت: یک الگوی ساختاری (router_prompt_template) برای هدایت مدل در دسته بندی پرسش ها به دسته های پیش فرض (VectorStore) ، SearchEngine یا None تعریف شد.

مدل زبان بزرگ: با استفاده از مدلی که در قدم (ب) ایجاد شده است، این بخش پرسش کاربر را بر اساس الگوی تعریف شده پردازش کرده و یک پاسخ دسته بندی تولید می کند.

پردازش خروجی: خروجی از مدل سپس با استفاده از Pydantic پردازش می شود تا مطمئن شود که با فرمت و دسته بندی مورد انتظار (QueryType class) سازگار است.

```

1 router_prompt_template = \
2     """
3     Your task is to categorize user queries into one of three categories: VectorStore, SearchEngine, or None.
4     Select VectorStore for questions related to Natural Language Processing or Speech Processing.
5     Choose SearchEngine for inquiries about computer science topics not involving NLP.
6     Opt for None if the question does not pertain to NLP or Computer Science.
7     Provide only the selected category as your response. Do not include any additional information.
8     {output_instruction}
9     query: {query}
10    """
11

```

```

1 router_prompt_template = dedent(router_prompt_template)
2
3 router_prompt = ChatPromptTemplate.from_template(
4     template=router_prompt_template
5 )

```

ت) برای مدیریت قابلیت‌های خروجی مدل، یک پارسر مبتنی بر Pydantic (router_parser) توسعه داده شد. این پارسر خروجی مدل را ارزیابی می‌کند تا آن را به یکی از دسته‌های مشخص شده (VectorStore)، SearchEngine یا None تبدیل کند.

```

1 class QueryType(BaseModel):
2     class_name: Literal["None", "SearchEngine", "VectorStore"] = Field()
3
4 router_parser = PydanticOutputParser(pydantic_object=QueryType)

```

```

1 # for test
2 router_prompt.invoke({
3     "query": "What is NLI?",
4     "output_instruction": router_parser.get_format_instructions()
5 })

```



```
1 router_chain = router_prompt | router_llm | router_parser
```

ث) تصمیم به تنظیم پارامتر دما به مقدار ۰/۰ برای این برنامه به دلایل استراتژیک اتخاذ شد. در زمینه استنتاج مدل، دمای ۰/۰ مدل را به سمت تولید پاسخ‌های قطعی و بدون متغیریت هدایت می‌کند به جای پاسخ‌های احتمالی.

این رفتار قطعی بسیار مهم است تا مدل با تحلیل دقیقی از پرسش‌های کاربر به یکی از دسته‌های مشخص شده (VectorStore)، SearchEngine یا None دسته‌بندی شود. این انتخاب روشی است که قابلیت‌های دسته‌بندی دقیق و قابل اعتماد چت‌بات ما را تقویت می‌کند، با تطابق که با نیازهای خاص پردازش زبان طبیعی و علوم کامپیوتر دارد.

با استفاده از این روش قطعی، مطمئن می‌شویم که مدل به طور ثابت یکی از دسته‌های مشخص شده (VectorStore)، SearchEngine یا None را بر اساس درک زبانی پرسش انتخاب می‌کند. این انتخاب روشی است که قابلیت‌های چت‌بات را در پاسخ به پرسش‌های متنوع کاربر بهبود می‌بخشد.

تست router_chain :



```
1 def test_router_chain(query):
2     return router_chain.invoke({
3         "query": query,
4         "output_instruction": router_parser.get_format_instructions()
5     })
6
7 print(
8     f"related_query: {test_router_chain('What are the main uses of Named Entity Recognizers?')}\n"
9     f"unrelated_query: {test_router_chain('Who are the top providers of cloud computing services?')}\n"
10    f"super_unrelated_query: {test_router_chain('Which teams have won the most Asian Cup titles?')}\n"
11 )
```

[22]

```
... related_query: class_name='VectorStore'
unrelated_query: class_name='SearchEngine'
super_unrelated_query: class_name='None'
```

بخش پنجم - پیاده‌سازی SEARCH ENGINE CHAIN

در این بخش، ما بر پیاده‌سازی زنجیره موتور جستجو با استفاده از پلتفرم Tavily تمرکز می‌کنیم که یک ابزار بازیابی مستندات آسان برای مدل‌های زبان بزرگ (LLMs) فراهم می‌کند. هدف این پیاده‌سازی، افزایش توانایی ما در پاسخ به پرسش‌های کاربران درباره علوم کامپیوتر، به ویژه آن‌هایی که خارج از حوزه پردازش زبان طبیعی (NLP) هستند، با بازیابی مستندات مرتبط از وب است.

الف) ابتدا باید در پلتفرم Tavily ثبت‌نام کنیم تا یک KeyAPI رایگان دریافت کنیم. این کلید برای تأیید درخواست‌های ما به موتور جستجوی Tavily استفاده خواهد شد.

```
1 os.environ["TAVILY_API_KEY"] = "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"
2 os.environ["TOGETHER_API_KEY"] = "9f24f1264668e13d395b5f214fce62f405544086d879b3fa954b5578de15b72d"
```

ب) پس از دریافت KeyAPI، ابزار جستجوی Tavily را در فضای Long Chain تعریف می‌کنیم. سپس یک پرسش نمونه به عنوان ورودی می‌دهیم تا خروجی را بررسی کنیم و با قالب داده‌های بازیابی شده آشنا شویم.

```
1 tavily_search_wrapper = TavilySearchAPIWrapper(tavily_api_key=os.environ["TAVILY_API_KEY"])
2 tavily_search = TavilySearchResults(api_wrapper=tavily_search_wrapper, max_results=5)
```

```
1 tavily_search.invoke("What are the latest advancements in NLP for improving sentiment analysis?")
```

```
[{'url': 'https://www.sciencedirect.com/science/article/pii/S2949719124000074',
  'content': 'Natural Language Processing Journal\nRecent advancements and challenges of NLP-based sentiment analysis: A state-of-the-art rev',
  'url': 'https://www.intechopen.com/chapters/87589',
  'content': 'The 2023 Sentiment Analysis Roadmap\nWritten By\nPublished: 10 January 2024\nDOI: 10.5772/intechopen.112276\nCite this chapter\
  'url': 'https://arxiv.org/pdf/2305.14842',
  'content': 'Žá\xa0EáYÉ\1a"ť"€"-Áq iUøIl\1fæD\1fÜ\Mž<½\noøÇ\1a0'F\1x13=\12øFÁp?øb"\"ty\"P\\xadÇ\1x13,ø,ø1"èð\11q\1x17pE:\nÇ aQgGDŸ\1x19\
  'url': 'https://www.nature.com/articles/s41598-024-60210-7',
  'content': 'Sentiment analysis, the computational task of determining the emotional tone within a text, has evolved as a critical subfield
  'url': 'https://www.researchgate.net/publication/378613766_Recent_advancements_and_challenges_of_NLP-based_sentiment_analysis_A_state-of-th
  'content': 'Sentiment analysis is a method within natural language processing that evaluates and identifies the emotional. tone or mood con
```

پ) وظیفه اصلی ما نوشتن زنجیره‌ای است که پرسش کاربر را به عنوان ورودی دریافت کرده و تا پنج مستند مرتبط را بازیابی کند. هر مستند بازیابی شده باید شامل محتوای مستند (page_content) و آدرس سایت مربوطه به عنوان metadata باشد.

مقدار حداکثر ۵ را در قسمت قبلی مشخص کردیم.



```
1 @chain
2 def parse_tavily_search(documents):
3     return [Document(
4         page_content=doc['content'],
5         metadata={'url': doc['url']}
6     ) for doc in documents]
7
8     return result_documents
```



```
1 search_engine_chain = tavily_search | parse_tavily_search
```



```
1 search_engine_chain.invoke("What are the latest advancements in NLP for improving sentiment analysis?")
```

```
[27] Python
... [Document(metadata={'url': 'https://www.sciencedirect.com/science/article/pii/S2949719124000074'}, page_content='Natural Language Processing
Document(metadata={'url': 'https://www.intechopen.com/chapters/87589'}, page_content='The 2023 Sentiment Analysis Roadmap\nWritten By\nPubli
Document(metadata={'url': 'https://www.startus-insights.com/innovators-guide/natural-language-processing-trends/'), page_content='Top 9 Natu
Document(metadata={'url': 'https://monkeylearn.com/blog/nlp-trends/'), page_content='December 23rd, 2020\nPosts you might like...\nNatural L
Document(metadata={'url': 'https://www.researchgate.net/profile/Jamin-Jim/publication/378613766_Recent_advancements_and_challenges_of_NLP-ba
```

بخش ششم - پیاده سازی RELEVANCY CHECK CHAIN (امتیازی)

در این بخش، ما بر روی پیاده سازی زنجیره ای تمرکز می کنیم که به منظور ارزیابی ارتباط اسناد بازاریابی شده با پرسش کاربر طراحی شده است. این زنجیره از یک مدل زبان بزرگ (LLM) استفاده می کند تا تعیین کند که آیا یک سند با پرسش کاربر مرتبط است یا خیر. این بررسی ارتباط برای بهبود کیفیت و قابلیت اطمینان پاسخ های تولید شده توسط چت بات بسیار مهم است.

الف) برای شروع، ما پرامپت هایی را طراحی می کنیم که به مدل زبان بزرگ دستور می دهد که چگونه ارتباط یک سند با پرسش کاربر را ارزیابی کند. این پرامپت به LLM اطلاع می دهد که وظیفه آن طبقه بندی سند به عنوان "مرتبط" یا "نامرتبط" است.

```

1 relevancy_check_template = \
2 """You are provided with a user question and a document.
3 If the given document is relevant to the user question and can be used to answer it,
4 output 'Relevant', and if not, output 'Irrelevant'. Only output the words Relevant and Irrelevant in a
5 JSON format as described in the output instructions.
6 User question: {user_query}
7 Document: {retrieved_document}
8 Output instruction: {output_instruction}
9 """
10

```

ب) ما از کتابخانه ChatTogether برای تعریف مدل زبان بزرگ استفاده می‌کنیم. مدل انتخاب شده meta-llama/Llama-3-70b-chat-hf است که به خاطر عملکرد قوی خود در وظایف درک زبان طبیعی شناخته شده است. پیکربندی مدل اطمینان می‌دهد که خروجی قطعی است با تنظیم دما به ۰.

```

1 relevancy_check_template = dedent(relevancy_check_template)
2 relevancy_check_prompt = ChatPromptTemplate.from_template(
3     template=relevancy_check_template
4 )
5
6 relevancy_check_llm = ChatTogether(
7     together_api_key=os.environ["TOGETHER_API_KEY"],
8     model="meta-llama/Llama-3-70b-chat-hf",
9     temperature=0
10 )

```

پ) خروجی از LLM نیاز به پردازش مناسب و تجزیه دارد. ما یک پردازشگر پس از مرتبط با استفاده از Pydantic تعریف می‌کنیم تا خروجی مدل را به یک فرمت ساختاریافته تبدیل کنیم.

```

1 class RelevancyType(BaseModel):
2     class_name: Literal["Relevant", "Irrelevant"] = Field()
3
4 relevancy_check_parser = PydanticOutputParser(pydantic_object=RelevancyType)

```

ت) اجرای زنجیره بررسی ارتباط برای اطمینان از اینکه پاسخ‌های ارائه شده توسط چت‌بات دقیق و مفید هستند ضروری است. برای مثال، اگر کاربری بپرسد: "استفاده‌های اصلی از شناساگرهای نام چیست؟"، و سند بازایی شده درباره باغبانی باشد، زنجیره این سند

را به درستی به عنوان "نامرتبط" طبقه‌بندی خواهد کرد. این کمک می‌کند تا کیفیت تعاملات حفظ شود و اطمینان حاصل شود که پرسش‌های کاربران به درستی مورد بررسی قرار می‌گیرند.

```
1 relevancy_check = relevancy_check_prompt | relevancy_check_llm | relevancy_check_parser
```

تست relevancy_check :

```
1 def test_relevancy_check(query, document=None):
2     if document == None:
3         document = search_engine_chain.invoke(query)[0]
4
5     result = relevancy_check.invoke({
6         "user_query": query,
7         "retrieved_document": document,
8         "output_instruction": relevancy_check_parser.get_format_instructions()
9     })
10
11     if result.class_name == 'Irrelevant':
12         print('Irrelevant query')
13         return result
14
15     print(f'Retrieved Document: {document}')
16     return result, document
```



```
> result, result_document = test_relevancy_check("What are the main uses of Named Entity Recognizers?")
result
[32]
... Retrieved Document: page_content='Unlock the Power of Visual Data Processing
Navigating the World of MLOps Certifications
Adel Nehme
10 min
How to Learn Machine Learning in 2024
Adel Nehme
15 min
Why ML Projects Fail, and How to Ensure Success with Eric Siegel, Founder of Machine Learning Week, Former
Adel Nehme
47 min
Multilayer Perceptrons in Machine Learning: A Comprehensive Guide
Sejal Jaiswal
15 min
A Beginner's Guide to CI/CD for Machine Learning
Abid Ali Awan
15 min
OpenCV Tutorial: Navigating the World of MLOps Certifications
How to Learn Machine Learning in 2024
Why ML Projects Fail, and How to Ensure Success with Eric Siegel, Founder of Machine Learning Week, Former
Multilayer Perceptrons in Machine Learning: A Comprehensive Guide
A Beginner's Guide to CI/CD for Machine Learning
OpenCV Tutorial: Clean the Text
In this section, we will clean our dataset using the NLTK library by following a few steps:
Entity Recognition
After adding a new pipeline to our model, we can visualize the named entities in our text using spaCy's d
...
Learn about how data is applied by industry leaders
Discover content by tools and technology

> test_relevancy_check("What are the best practices for growing tomatoes in a home garden?", document=result_document)
[33]
... Irrelevant query
... RelevancyType(class_name='Irrelevant')
```

بخش هفتم - پیاده سازی FALLBACK CHAIN

در این بخش، هدف ما پیاده سازی یک مکانیزم پشتیبان برای چت بات است تا به طور مناسب به پرسش هایی که خارج از حوزه تخصصی آن هستند، پاسخ دهد. مسیر یاب ما سه خروجی ممکن دارد: ذخیره سازی برداری، استفاده از موتور جستجو به عنوان پشتیبان، و عدم پاسخگویی. ما قبلاً دو مورد اول را پیاده سازی کرده ایم و اکنون بر روی گزینه ی آخر تمرکز می کنیم. زنجیره پشتیبان به چت بات اجازه می دهد تا با احترام از پاسخ دادن به سوالاتی که خارج از دامنه تخصصی آن هستند، خودداری کند.

الف) ما نیاز داریم که یک الگوی پرسش ایجاد کنیم که حوزه دانش چت بات را به وضوح تعریف کند. این الگو شامل مکان نگارهایی برای تاریخچه گفتگو و پرسش کاربر است. این الگو به مدل دستوری می دهد که به شیوه ای مودبانه توجه کند که نمی تواند به پرسش هایی که خارج از حوزه NLP یا تشخیص گفتار هستند، پاسخ دهد.

```

1 fallback_template = """You are a helpful and knowledgeable teaching assistant.
2 Your role is to supply educational information focused on Natural Language Processing (NLP)
3 and Speech Recognition to the human user. Avoid addressing questions that fall outside the
4 scope of NLP and Speech Recognition. If a query is not relevant, please acknowledge this limitation.
5 Current conversation:
6 {chat_history}
7 Human: {query}
8 """

```

ب) سپس، ما یک مدل زبان با استفاده از تنظیمات "Together" تعریف می‌کنیم و الگوی پرسش را به آن متصل می‌کنیم. ما از مقدار دمای بالاتر استفاده می‌کنیم تا به مدل این امکان را بدهیم که پاسخ‌های متنوع و خلاقانه‌تری تولید کند.

```

1 @chain
2 def chat_history_gather(state):
3     return {
4         "chat_history": gather_chat_history(state),
5         "query": state['query']
6     }
7
8 def gather_chat_history(state):
9     return [(
10         f"human: {msg.content}"
11         if isinstance(msg, HumanMessage)
12         else f"AI: {msg.content}"
13     ) for msg in state["chat_history"]]
14
15 fallback_template = dedent(fallback_template)
16 fallback_prompt = ChatPromptTemplate.from_template(
17     fallback_template
18 )
19
20 fallback_llm = ChatTogether(
21     together_api_key=os.environ["TOGETHER_API_KEY"],
22     model="meta-llama/Llama-3-70b-chat-hf",
23     temperature=1
24 )

```

پ) در نهایت، خروجی مدل را به یک پارسر خروجی متنی متصل می‌کنیم که پاسخ مدل را به یک فرمت مهره‌ای برای کاربر تبدیل می‌کند.

```
1 fallback_chain = chat_history_gather | fallback_prompt | fallback_llm | StrOutputParser()
```

ما همچنین نیاز داریم که یک دیکشنری سفارشی برای ذخیره تاریخچه چت و موارد دیگر مربوط به متن محیط گفتگو ایجاد کنیم.

```
1 class AgentState(TypedDict):
2     """Represents the current chat context and history with the agent."""
3     query: str
4     generation: str
5     documents: list[Document]
6     chat_history: list[BaseMessage]
```

در زیر مثالی از استفاده از زنجیره پشتیبان با یک پرسش که خارج از حوزه NLP و تشخیص گفتار است، آورده شده است.

```
1 state = AgentState(
2     query="What are the best practices for growing tomatoes in a home garden?",
3     generation="",
4     documents=[],
5     chat_history=[]
6 )
7
8 fallback_chain.invoke(state)
```

زمانی که زنجیره پشتیبان با یک پرسش که خارج از حوزه NLP و تشخیص گفتار است، فراخوانی می‌شود، مدل یک پاسخ تولید می‌کند که محدودیت خود را به آن اعلام کرده و پیشنهاد می‌دهد که در موضوعات مربوط به NLP یا تشخیص گفتار، پاسخ دهد. به عنوان مثال:

```
[38]
... "I'm happy to help you with your question! However, I must acknowledge that my expertise lies in Natural Language Processi
```

این روش تضمین می‌کند که چت‌بات در حوزه تخصصی خود باقی می‌ماند و در مواجهه با پرسش‌هایی که خارج از دامنه دانش خود هستند، پاسخ مودبانه و اطلاع‌رسانی ارائه می‌دهد.

بخش هشتم - پیاده سازی GENERATIVE WITH CONTEXT CHAIN

در این بخش، ما قصد داریم یک زنجیره چت بات به نام `generate_with_context_chain` پیاده سازی کنیم. هدف از این زنجیره این است که سوال کاربر را به همراه اسناد مرتبط به عنوان ورودی بگیرد و پاسخ نهایی را ارائه دهد. این برای توسعه یک دستیار هوشمند و آگاه به زمینه که می تواند پاسخ های دقیقی بر اساس اطلاعات ارائه شده تولید کند، بسیار مهم است.

الف) پرامپت این زنجیره شامل ارائه سوال کاربر و مجموعه ای از پاسخ های ممکن در زمینه اسناد ارائه شده به مدل است. وظیفه مدل این است که با استفاده از اطلاعات موجود در اسناد، یک پاسخ دقیق تولید کند. اگر اطلاعات لازم در زمینه ارائه شده موجود نباشد، مدل باید کاربر را به درستی مطلع کرده و از ارائه هرگونه اطلاعات خارج از زمینه خودداری کند.

```
1 generate_template = """You are an intelligent and helpful assistant.
2 Answer the user's query using only the information provided in the
3 given context. If the context does not contain the necessary information,
4 inform the user that you cannot answer based on the available context,
5 and refrain from providing information outside of it.
6 Context: {documents}
7 Query: {query}
8 """
```

ب) مدل استفاده شده در این زنجیره با مدل های استفاده شده در بخش های قبلی یکسان است. این تضمین می کند که عملکرد و یکپارچگی سیستم به صورت یکنواخت حفظ شود.

```
1 generate_template = dedent(generate_template)
2 generate_prompt = ChatPromptTemplate.from_template(generate_template)
```

پ) `StrOutputParser` به عنوان پس پردازنده برای این زنجیره استفاده می شود. این پردازشگر خروجی مدل را به فرمت رشته ای کاربرپسند تبدیل می کند.

```
1 generate_chain = generate_prompt | fallback_llm | StrOutputParser()
```

تست برنامه:

```
state = AgentState(
... query="What are the best practices for growing tomatoes in a home garden?",
... generation="",
... documents=[result_document],
... chat_history=[]
)

generate_chain.invoke(state)
```

[42]

```
... 'I cannot answer your query based on the available context. The provided context appears to be re
```

بخش نهم – آماده سازی گراف با استفاده از LANGGRAPH

در این بخش، گراف یک عامل گفتگو را با استفاده از ابزار Line Graph می‌سازیم و توضیح می‌دهیم. این گراف برای مدیریت و پردازش پرسش‌های کاربران از طریق گره‌های مختلف طراحی شده است که هر کدام مسئول انجام وظایف خاصی در جریان مکالمه هستند.

جریان کاری گراف

جریان کاری با پردازش پرسش کاربر توسط گره Router آغاز می‌شود که پرسش را به یکی از سه مسیر هدایت می‌کند Vector : Search Engine یا Fallback. بسته به مسیری که انتخاب شده، پرسش از طریق گره‌های بعدی پردازش می‌شود تا پاسخ نهایی تولید و به کاربر بازگردانده شود. فرآیند زمانی پایان می‌یابد که پاسخ نهایی تولید و به کاربر ارائه شود.

پیاده‌سازی گره‌ها

```
1 def router_node(state: dict):
2     try:
3         response = question_router.invoke({"query": query,
4                                             "output_instructions": question_router_parser.get_format_instructions()
5                                             })
6     except:
7         return 'LLMFallback'
8
9     return 'LLMFallback' if response.class_name == 'None' else response.class_name
```

ایجاد گراف عامل

```

1 def filter_documents(state: dict):
2     relevant_documents = []
3     for document in state['documents']:
4         try:
5             relevancy_result = relevancy_check.invoke({
6                 'user_query': state['query'],
7                 'retrieved_document': document,
8                 'output_instruction': relevancy_check_parser.get_format_instructions()
9             })
10            if relevancy_result.class_name == 'Relevant':
11                relevant_documents.append(document)
12        except Exception:
13            continue
14
15    return relevant_documents

```

```

1 workflow = StateGraph(AgentState)
2
3 ensemble_retriever = build_ensemble_retriever([0.5, 0.5])
4 workflow.add_node('vector_store', lambda state: {'documents': ensemble_retriever.invoke(input=state['query'])})
5 workflow.add_node('search_engine', lambda state: {'documents': search_engine_chain.invoke(state['query'])})
6 workflow.add_node('fallback', lambda state: {'generation': fallback_chain.invoke(state)})
7 workflow.add_node('generate_with_context', lambda state: {'generation': generate_chain.invoke(state)})
8 workflow.add_node('filter_docs', lambda state: {'documents': filter_docs_node(state)})
9
10 workflow.set_conditional_entry_point(
11     router_node,
12     {
13         'VectorStore': 'vector_store',
14         'SearchEngine': 'search_engine',
15         'LLMFallback': 'fallback',
16     }
17 )
18
19 workflow.add_edge('vector_store', 'filter_docs')
20 workflow.add_edge('search_engine', 'filter_docs')
21 workflow.add_edge('fallback', END)
22 workflow.add_edge('generate_with_context', END)
23
24 workflow.add_conditional_edges(
25     'filter_docs',
26     lambda docs: 'generate_with_context' if len(docs) != 0 else 'search_engine',
27     {
28         'search_engine': 'search_engine',
29         'generate_with_context': 'generate_with_context'
30     }
31 )

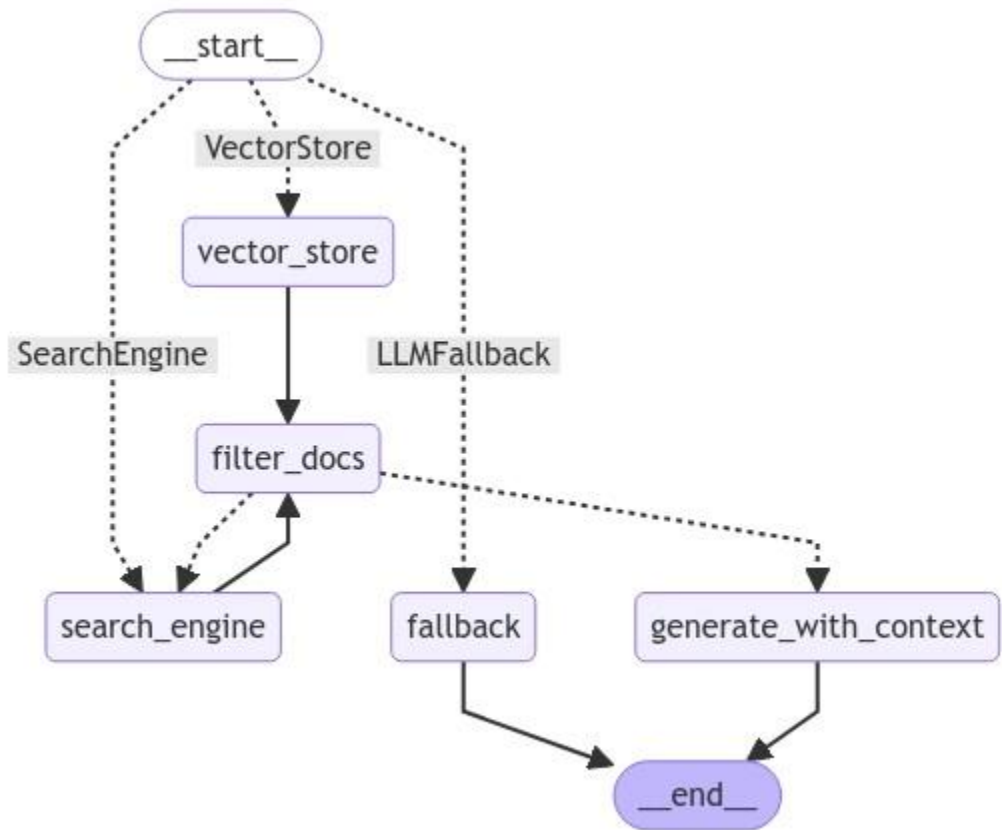
```

نمایش گراف

```

1 app = workflow.compile(debug=False)
2 display(Image(app.get_graph().draw_mermaid_png(draw_method=MermaidDrawMethod.API)))

```



آزمایش گراف و گفتگو با چت بات

```

1 def send_new_message(query, state=None):
2     if state == None:
3         state = {}
4         state['chat_history'] = []
5
6     response = app.invoke({
7         'query': query,
8         'chat_history': state['chat_history']
9     })
10
11     response['chat_history'] = [HumanMessage(response['query']), AIMessage(response['generation'])]
12     return response
  
```

```
state = send_new_message("Hello Atask! I want to know about LSTM")
Markdown(state['generation'])
```

[48]

... Hello!

LSTM (Long Short-Term Memory) is a fascinating topic in the realm of Natural Language Processing (NLP). LSTM is a type of Recurrent Neural Network (RNN) architecture that's particularly well-suited for modeling sequential data. The main advantage of LSTM over traditional RNNs is its ability to learn long-term dependencies in data. This is achieved through the use of memory cells, which allow the network to store information for extended periods of context and relationships between words or tokens can span multiple time steps.

In the context of speech recognition, LSTM is often used as a key component in acoustic models, which aim to predict phonemes or speech sounds from audio inputs. LSTM's ability to capture long-term dependencies helps to improve the accuracy of speech recognition systems.

Some of the key benefits of LSTM include:

1. **Improved performance:** LSTM outperforms traditional RNNs in many NLP tasks, particularly those involving long-range dependencies.
2. **Handling vanishing gradients:** LSTM's memory cells help to mitigate the vanishing gradient problem, which occurs when gradients are backpropagated through time, causing them to shrink exponentially.
3. **Modeling complex patterns:** LSTM is capable of modeling complex patterns and relationships in sequential data, making it a powerful tool for tasks like language modeling, machine translation, and text classification.

Would you like to know more about how LSTM is used in specific NLP tasks, such as language modeling or machine translation?

```
state = send_new_message("Now I want to know about forgive gate", state)
Markdown(state['generation'])
```

[49]

... The "forget gate" is a crucial component of the LSTM architecture.

In an LSTM network, the forget gate is one of the three gates that regulate the flow of information into and out of the memory cell. The forget gate is responsible for deciding what information to discard from the cell.

The forget gate is a sigmoid neural network layer that takes the previous hidden state and the current input as inputs. It outputs a vector of values between 0 and 1, which are then element-wise multiplied with the previous memory cell state, effectively "forgetting" or removing some information from the memory cell.

The forget gate plays a crucial role in learning long-term dependencies, as it allows the LSTM network to selectively retain or discard information that is no longer relevant and reduce the impact of irrelevant information.

In the context of speech recognition, the forget gate helps the LSTM network to forget some of the acoustic features that are not relevant for predicting the next word in the sequence.

Would you like to know more about the other gates in the LSTM architecture, such as the input gate and the output gate?

```
state = send_new_message("Now tell me about new mechanisms for gates in LSTM", state)
Markdown(state['generation'])
```

[50]

... In recent years, several new mechanisms have been proposed to enhance the functionality of gates in LSTM networks. Here are a few examples:

1. **Gated Recurrent Units (GRUs):** While not exactly a modification to the LSTM gate, GRUs are a related type of RNN that use gates in a slightly different way. GRUs have two gates: a reset gate to determine what information to forget, while the update gate determines how much new information to add. GRUs are simpler and faster to train than LSTMs but still capture long-term dependencies effectively.
2. **Coupled Input and Forget Gates:** In traditional LSTMs, the input and forget gates are learned independently. However, some researchers have proposed coupling these gates, so that the network can better balance the amount of information to forget and the amount of new information to add.
3. **Gate Attention Mechanisms:** This mechanism involves adding attention weights to the gates, allowing the network to focus on specific parts of the input sequence when deciding what to forget or what to add. This is particularly useful in tasks like machine translation, where certain acoustic features may be more relevant than others.
4. **Learned Normalize Gate:** This mechanism involves adding a learnable normalization factor to the gates, allowing the network to adjust the scale of the gate outputs. This can help to stabilize the training process and improve the network's performance on long-term dependencies.
5. **Phase-LSTM:** This is a type of LSTM that uses a phase component to modulate the gates. The phase component is a learnable parameter that captures the periodic patterns in the input sequence, allowing the network to better model the temporal structure of the data.
6. **Gating Mechanisms with External Memory:** This involves using external memory to store information that can be used to inform the gating mechanism. For example, the network could store information about the current state of the input sequence, and then use this information to decide what to remember or forget.

These are just a few examples of the new mechanisms that have been proposed to enhance the functionality of gates in LSTM networks. Would you like to know more about any of these mechanisms?

پایان!