# Simple_RAG 实现

## 1. 配置虚拟环境，安装依赖

### 1.1 构建环境

```
1 conda create -n rag_test python=3.10
2 source activate rag_test
```

### 1.2 安装依赖

requirements.txt 文件

```
1 faiss-cpu
2 faiss-gpu==1.7.2
3 argparse==1.1
4 transformers==4.37.2
5 torch==2.0.0
6 tqdm==4.66.2
```

运行以下脚本

```
1 pip install -r requirements.txt -i https://pypi.tuna.tsinghua.edu.cn/simple --ignore-installed
```

## Embedding 模型

- embedding模型使用下载好的m3e-large模型

```
1 tokenizer = AutoTokenizer.from_pretrained('m3e-large')
2 model = AutoModel.from_pretrained('m3e-large').to('cuda:3')
```

# 2. 代码

## 2.1 导包

```
1  import os
2  import argparse
3  import faiss
4  from transformers import AutoTokenizer, AutoModel, AutoModelForCausalLM
5  import torch
6  from tqdm import tqdm
```

### 2.1.1 加载transformers出现报错

```
1  ImportError: cannot import name 'COMMON_SAFE_ASCII_CHARACTERS' from
   'charset_normalizer.constant'
```

- 解决方案

```
1  pip install --upgrade charset-normalizer
```

## 2.2 文件处理函数

```
1  def process_file(file_path):
2      with open(file_path, encoding = 'utf-8') as f:
3          text = f.read()
4          sentences = text.split('\n')
5          return text, sentences
6
7  process_file(os.path.join(os.getcwd(),'requirements.txt'))
```

Output:

('faiss-cpu\nfaiss-gpu==1.7.2\nargparse==1.1\ntransformers==4.37.2\ntorch==2.0.0\ntqdm==4.66.2',
 ['faiss-cpu',
  'faiss-gpu==1.7.2',
  'argparse==1.1',
  'transformers==4.37.2',
  'torch==2.0.0',
  'tqdm==4.66.2'])

## 2.3 构建prompt函数

```python
def generate_rag_prompt(data_point):
    return f'''### Instruction:
{data_point['instruction']}
### Input:
{data_point['input']}
### Response:
    '''
```

## 2.4 生成embedding

```python
class DocumentEmbedder:
    def __init__(self, max_length = 128, max_number_of_sentences = 3):
        self.model = model
        self.tokenizer = tokenizer
        self.max_length = max_length
        self.max_number_of_sentences = max_number_of_sentences

    def get_document_embeddings(self,sentences):
        sentences = sentences[:self.max_number_of_sentences]
        encoded_input = self.tokenizer(sentences,
                                        padding = True,
                                        truncation = True,
                                        max_length = 128,
                                        return_tensors = 'pt')
        encoded_input.to('cuda:3')
        with torch.no_grad():
            model_output = self.model(**encoded_input)
        return torch.mean(model_output.pooler_output,dim = 0,keepdim = True)
```

## 2.5 提前准备好知识库

```python
1  from langchain_text_splitters import RecursiveCharacterTextSplitter
2  from tqdm import tqdm
3  text_splitter = RecursiveCharacterTextSplitter(
4      # Set a really small chunk size, just to show.
5      chunk_size=500,
6      chunk_overlap=20,
7      length_function=len,
8      is_separator_regex=False,
9  )
10 texts = text_splitter.split_text(text)
11 for i in range(len(texts)):
12     with open(f'RAG-file/chunks/chunks-{i}','w') as f:
13         f.write(texts[i])
```

- 将一篇长文本txt文件进行500个词分词写入chunks文件夹内，本次使用的文本是self-RAG那篇文章的文本

## 2.6 大模型类

- 这里调用了zhipuAI的API

# 主函数定义

```python
1  if __name__ == '__main__':
2      documents = {}
3      for idx, file in enumerate(tqdm(os.listdir('RAG-file/chunks'))):
4          current_filepath = os.path.join('RAG-file/chunks',file)
5          text, sentences = process_file(current_filepath)
6          documents[idx] = {'file_path':file,
7                            'sentences':sentences,
8                            'document_text':text}
9
10     print('Getting document embeddings...')
11     document_embedder = DocumentEmbedder()
12     embeddings = []
13     for idx in tqdm(documents):
14
   embeddings.append(document_embedder.get_document_embeddings(documents[idx]
   ['sentences']))
15     embeddings = torch.concat(embeddings,dim = 0).data.cpu().numpy()
```

```
16        embedding_dimensions = embeddings.shape[1]
17        faiss_index = faiss.IndexFlatIP(int(embedding_dimensions))
18        faiss_index.add(embeddings)
19
20        question = 'Introduce to me what is RAG and how to use it'
21
22        query_embedding = document_embedder.get_document_embeddings([question])
23        distances, indices =
      faiss_index.search(query_embedding.data.cpu().numpy(),k=int(5))
24
25        context = ''
26        for idx in indices[0]:
27            context += documents[idx]['document_text']
28
29        rag_prompt = generate_rag_prompt({'instruction':question,
30                                          'input':context})
31        print('Generating answer...')
32        answer = GenerativeModel(rag_prompt)
33        print(answer)
```

# 输出结果

- 输入问题:

- Introduce to me what is RAG and how to use it

- 输出:

```
100%|██████████| 170/170 [00:00<00:00, 60472.54it/s]
Getting document embeddings...
100%|██████████| 170/170 [00:03<00:00, 50.69it/s]
Generating answer...
RAG, or Retrieval-augmented Generation, is a technique that combines the power of large language models (LLMs) with a retrieval component to enhance the generation p
rocess. The retrieval component allows the model to access relevant information from a large database or corpus, which can then be used to improve the coherence, acc
uracy, and diversity of the generated text.

The basic idea behind RAG is to use the LLM to generate text while periodically consulting the retrieval component for additional context or information. This can he
lp the model to generate more accurate and relevant outputs, especially in tasks that require a deep understanding of specific domains or a large amount of context.

To use RAG, you will need to have access to a large language model and a retrieval component. The retrieval component can be a simple database, a pre-trained knowled
ge graph, or even another model that has been trained to represent knowledge in a useful way.

Once you have access to these components, you can use the RAG technique by having the LLM generate text while periodically providing it with additional context or in
formation from the retrieval component. This can be done by either having the LLM query the retrieval component directly, or by having a human provide the additional
context.

RAG has shown promise in a variety of tasks, including machine translation, summarization, and question-answering. It is particularly useful in tasks where the model
needs to access a large amount of context or domain-specific knowledge to generate accurate and relevant outputs.
```