

基于Zero-shot方式实现LLM文本匹配

学习目标

- 掌握Zero-shot方式下prompt的设计方式
- 掌握利用LLM实现文本匹配的代码

1 LLM信息抽取任务介绍

- 首先，我们构造几个短文本对：

```
1 1. ('股票市场今日大涨，投资者乐观。', '持续上涨的市场让投资者感到满意。'),  
2 2. ('油价大幅下跌，能源公司面临挑战。', '未来智能城市的建设趋势愈发明显。'),  
3 3. ('利率上升，影响房地产市场。', '高利率对房地产有一定冲击。'),
```

- 我们期望模型能够帮我们识别出这 3 对句子中，哪几对描述的是相似的语言。
- 我们期望模型输出的结果为：

```
1 ['相似', '不相似', '相似']
```

2 Prompt设计

- 在该任务的 prompt 设计中，我们主要考虑 2 点：
 - 需要向模型解释什么叫作「文本匹配任务」
 - 需要让模型按照我们指定的格式输出
- 为了让模型知道什么叫做「文本匹配任务」，我们借用 Incontext Learning 的方式，先给模型展示几个正确的例子：

```
1 >>> User: 句子一：公司ABC发布了季度财报，显示盈利增长。\\n句子二：财报披露，公司ABC利润上升  
2 >>> Bot: 是  
3 >>> User: 句子一：黄金价格下跌，投资者抛售。\\n句子二：外汇市场交易额创下新高  
4 >>> Bot: 不是  
5 ...
```

其中，`User` 代表我们输入给模型的句子，`Bot` 代表模型的回复内容。

注意：上述例子中 `Bot` 的部分也是由人工输入的，其目的是希望看到在看到类似 `User` 中的句子时，模型应当做出类似 `Bot` 的回答。

3 文本匹配任务代码实现

- 本章节使用的模型为ChatGLM-6B，参数参数较大（6B），下载到本地大概需要 12G+ 的磁盘空间，请确保磁盘有充足的空间。此外，加载模型大概需要 13G 左右的显存，如果您显存不够，可以进行模型量化加载以缩小模型成本。
- 本次文本匹配任务实现的主要过程：

- 构造prompt
- 实现文本匹配
- 代码存放位置: /Users/***/PycharmProjects/llm/zero-shot/llm_text_matching.py
- llm_information_extraction.py脚本中包含三个函数: init_prompts()和inference()

3.1 导入必备的工具包

```
1 from rich import print
2 from transformers import AutoTokenizer, AutoModel
3
4 import os
5
6
7 # 提供相似, 不相似的语义匹配例子
8 examples = {
9     '是': [
10         ('公司ABC发布了季度财报, 显示盈利增长。', '财报披露, 公司ABC利润上升。'),
11     ],
12     '不是': [
13         ('黄金价格下跌, 投资者抛售。', '外汇市场交易额创下新高。'),
14         ('央行降息, 刺激经济增长。', '新能源技术的创新。')
15     ]
16 }
17
```

3.2 构建init_prompts()函数

- 目的: 进行prompt设计
- 具体代码实现:

```
1 def init_prompts():
2     """
3     初始化前置prompt, 便于模型做 incontext learning.
4     """
5     pre_history = [
6         (
7             '现在你需要帮助我完成文本匹配任务, 当我给你两个句子时, 你需要回答我这两句话语义是否相似。'
8             '只需要回答是否相似, 不要做多余的回答。',
9             '好的, 我将只回答“是”或“不是”。'
10        )
11    ]
12    for key, sentence_pairs in examples.items():
13        for sentence_pair in sentence_pairs:
14            sentence1, sentence2 = sentence_pair
15            pre_history.append((
16                f'句子一: {sentence1}\n句子二: {sentence2}\n上面两句话是相似的语义吗? ',
17                key
18            ))
```

```
19
20     return {'pre_history': pre_history}
```

3.3 构建inference()函数

- 目的：模型实现信息匹配
- 具体代码实现

```
1  def inference(
2      sentence_pairs: list,
3      custom_settings: dict
4  ):
5      """
6      推理函数。
7
8      Args:
9          model (transformers.AutoModel): Language Model 模型。
10         sentence_pairs (List[str]): 待推理的句子对。
11         custom_settings (dict): 初始设定, 包含人为给定的 few-shot example.
12     """
13     for sentence_pair in sentence_pairs:
14         sentence1, sentence2 = sentence_pair
15         sentence_with_prompt = f'句子一: {sentence1}\n句子二: {sentence2}\n上面两句话是相似
16         的语义吗? '
17         response, history = model.chat(tokenizer, sentence_with_prompt,
18         history=custom_settings['pre_history'])
19         print(f'>>> [bold bright_red]sentence: {sentence_pair}')
20         print(f'>>> [bold bright_green]inference answer: {response}')
21         # print(history)
```

- 代码调用

```
1  if __name__ == '__main__':
2      #device = 'cuda:0'
3      device = 'cpu'
4      tokenizer = AutoTokenizer.from_pretrained("./ChatGLM-6B/THUDM/chatglm-6b",
5      trust_remote_code=True)
6      #model = AutoModel.from_pretrained("./ChatGLM-6B/THUDM/chatglm-6b",
7      # trust_remote_code=True).half().cuda()
8      model = AutoModel.from_pretrained("./ChatGLM-6B/THUDM/chatglm-6b",
9      trust_remote_code=True).float()
10     model.to(device)
11
12     sentence_pairs = [
13         ('股票市场今日大涨, 投资者乐观.', '持续上涨的市场让投资者感到满意.'),
14         ('油价大幅下跌, 能源公司面临挑战.', '未来智能城市的建设趋势愈发明显.'),
15         ('利率上升, 影响房地产市场.', '高利率对房地产有一定冲击.'),
16     ]
17
18     custom_settings = init_prompts()
19     inference(
20         sentence_pairs,
```

```
21 | custom_settings
22 | )
```

- 打印结果:

```
Leading checkpoint shards: 100% | 8/8 [00:12<00:00, 1.12s/it]
The dtype of attention mask (torch.int64) is not bool
>>> sentence: ('股票市场今日大涨，投资者乐观。', '持续上涨的市场让投资者感到满意。')
>>> inference answer: 是
>>> sentence: ('油价大幅下跌，能源公司面临挑战。', '未来智能城市的建设趋势愈发明显。')
>>> inference answer: 不是
>>> sentence: ('利率上升，影响房地产市场。', '高利率对房地产有一定冲击。')
>>> inference answer: 是
```

小结总结

本章节主要介绍了如何利用Few-Shot方式基于ChatGLM-6B实现文本匹配任务。