

天池SMP 2023 ChatGLM金融大模型挑战赛

馒头科技团队方案代码 B/C榜 TOP1

主要代码说明

- main.py
 - 包含全流程, 数据下载/pdf表格提取/表格生成/问题分类/问题关键词提取/SQL生成/问答结果生成/提交生成
- generate_answer_with_classify.py
 - generate_answer函数对每类问题生成问答结果

主要思路

问题分类

SQL生成

问答结果生成

Type1

核心处理代码

```
1 table_dict = {
2     'A': ['basic_info'],
3     'B': ['employee_info', 'dev_info'],
4     'C': ['cbs_info', 'cscf_info', 'cis_info'],
5     'G': ['basic_info', 'employee_info', 'dev_info', 'cbs_info',
6         'cscf_info', 'cis_info']
7 }
8 if len(company_abbrs) == 0:
9     logger.warning('匹配到了类别{}, 但是不存在报表'.format(question_type))
10 else:
11     # _, question_keywords = question_util.parse_question_keywords(model,
12     ori_question, real_comp, years)
13     logger.info('问题关键词: {}'.format(question_keywords))
14
```

```

14     background = ''
15     tot_matched_rows = []
16     for year in years:
17         pdf_table = load_tables_of_years(company, [year], pdf_tables,
18 pdf_info)
19
20         background += '已知{}(简称:{}, 证券代码: {}){}年的资料如下:\n
21 '.format(company, abbr, code, year)
22         matched_table_rows = []
23         for keyword in question_keywords:
24             matched_table_rows.extend(recall_pdf_tables(keyword, [year],
25 pdf_table,
26             min_match_number=3,
27 valid_tables=table_dict[question_type]))
28
29         if len(matched_table_rows) == 0:
30             for table_row in pdf_table:
31                 if table_row[0] in table_dict[question_type]:
32                     matched_table_rows.append(table_row)
33
34         table_text = table_to_text(real_comp, ori_question,
35 matched_table_rows, with_year=False)
36         background += table_text
37         background += '\n'
38
39         tot_matched_rows.extend(matched_table_rows)
40
41         tot_matched_rows = add_text_compare_in_table(tot_matched_rows)
42         tot_text = table_to_text(real_comp, ori_question, tot_matched_rows,
43 with_year=True)
44
45         if '相同' in tot_text or '不相同且不同' in tot_text:
46             answer = tot_text
47         else:
48             question_for_model = type1.get_prompt(ori_question, company, abbr,
49 years).format(background, ori_question)
50             logger.info('Prompt length {}'.format(len(question_for_model)))
51             if len(question_for_model) > 5120:
52                 question_for_model = question_for_model[:5120]
53             logger.info(question_for_model.replace('<', ''))
54             answer = model(question_for_model)

```

主要处理步骤:

1. 对于Type1每类问题, 从该问题对应的报表中根据关键词召回可能对应的行

2. 对于行为字符串类型的, add_text_compare_in_table实现对比不同年份的字段内容是否相同
3. 对于问题问是否相同的, 直接返回召回的报表
4. 对于其他的, 通过prompt组合召回的报表来传给模型进行回答, 具体prompt函数见 type1.get_prompt

Type2

核心处理代码

```
1 if type2.is_type2_growth_rate(ori_question):
2     years_of_table = []
3     for year in years:
4         years_of_table.extend([year, str(int(year)-1)])
5     pdf_table = load_tables_of_years(company, years_of_table,
6 pdf_tables, pdf_info)
7     pdf_table = add_growth_rate_in_table(pdf_table)
8 elif type2.is_type2_formula(ori_question):
9     pdf_table = load_tables_of_years(company, years, pdf_tables,
10 pdf_info)
11 else:
12     logger.error('无法匹配, 该问题既不是增长率也不是公式计算')
13     pdf_table = load_tables_of_years(company, years, pdf_tables,
14 pdf_info)
15
16 step_questions, step_keywords, variable_names, step_years, formula,
17 question_formula = type2.get_step_questions(
18     ori_question, ''.join(question_keywords), real_comp, years[0])
19 step_answers = []
20 variable_values = []
21 if len(step_questions) > 0:
22     for step_question, step_keyword, step_year in zip(step_questions,
23 step_keywords, step_years):
24         if len(step_keyword) == 0:
25             logger.error('关键词为空')
26
27         background = '已知{}{}年的资料如下:\n'.format(real_comp,
28 step_year)
29         # background += '-----\n'
30
31         matched_table_rows = recall_pdf_tables(step_keyword,
32 [step_year], pdf_table,
33         min_match_number=3, top_k=5)
34         # print(matched_table_rows)
```

```

28         if len(matched_table_rows) == 0:
29             logger.warning('无法匹配keyword {}, 尝试不设置限
30 制'.format(step_keyword))
31             matched_table_rows = recall_pdf_tables(step_keyword,
32 [step_year], pdf_table,
33             min_match_number=2, top_k=None)
34         if len(matched_table_rows) == 0:
35             logger.error('仍然无法匹配keyword {}'.format(step_keyword))
36             matched_table_rows = recall_pdf_tables(step_keyword,
37 [step_year], pdf_table,
38             min_match_number=0, top_k=10)
39
40         table_text = table_to_text(real_comp, ori_question,
41 matched_table_rows, with_year=False)
42         if table_text != '':
43             background += table_text
44
45         question_for_model =
46         prompt_util.get_prompt_single_question(ori_question, real_comp,
47 step_year).format(background, step_question)
48         logger.opt(colors=True).info('<cyan>{}
49 </>'.format(question_for_model.replace('<', '')))
50         step_answer = model(question_for_model)
51         variable_value =
52         type2.get_variable_value_from_answer(step_answer)
53         if variable_value is not None:
54             step_answers.append(step_answer)
55             variable_values.append(variable_value)
56             logger.opt(colors=True).info('<green>{}</><red>{}
57 </>'.format(step_answer.replace('<', ''), variable_value))
58         if len(step_questions) == len(variable_values):
59             for name, value in zip(variable_names, variable_values):
60                 formula = formula.replace(name, value)
61             result = None
62             try:
63                 result = eval(formula)
64             except:
65                 logger.error('Eval formula {} failed'.format(formula))
66             if result is not None:
67                 answer = ''.join(step_answers)
68                 answer += question_formula
69                 answer += '得出结果{:.2f}({:.2f}%)'.format(result, result*100)

```

主要处理步骤:

1. 如果是计算增长率, 则需要召回前一年的数据
2. type2.get_step_questions将需要计算的内容转换为多个提问, 得到计算公式中的每个元素, 传给模型进行回答后, 提取数值作为公式的单元, 最后通过python计算该公式
 - a. 例如对于增长率, 公式为 $(A-B)/B$, A表示当年的数值, B表示上年的数值
 - b. 通过prompt_util.get_prompt_single_question生成针对A和B的提问
 - c. 提取回答中的数值type2.get_variable_value_from_answer
 - d. eval(formula)得到计算的结果

Type3

✓ 核心代码

```
1 any_question, _ = question_util.parse_question_keywords(model,
2   ori_question, real_comp, years)
3
4 background = '*****{}{}年年报*****\n'.format(
5   real_comp, years[0])
6 matched_text = recall_annual_report_texts(model, any_question,
7   ''.join(question_keywords),
8   matched_pdf_names[0], None)
9 for block_idx, text_block in enumerate(matched_text):
10     background += '{}片段: {}{}\n'.format('-'*15, block_idx+1, '-'*15)
11     background += text_block
12     background += '\n'
13 question_for_model = prompt_util.prompt_question_tp31.format(
14     background, ori_question, ''.join(question_keywords),
15     ''.join(question_keywords), ''.join(question_keywords))
16
17 logger.info('Prompt length {}'.format(len(question_for_model)))
18 if len(question_for_model) > 5120:
19     question_for_model = question_for_model[:5120]
20 logger.info(question_for_model.replace('<', ''))
21 answer = model(question_for_model)
```

主要处理步骤:

1. question_util.parse_question_keywords提取问题中的关键词
2. recall_annual_report_texts进行文本的召回
 - a. 其中召回的主要代码如下, 这里按行读取pdf文件中的文本内容, 采用bm25算法, 分别对问题关键词和原问题进行召回, 然后合并召回的文本块, 最后按照和原问题的字符匹配长度取匹配

度最高的文本块:

▼

```
1 text_pages = load_pdf_pages(key)
2 text_lines = list(itertools.chain(*[page.split('\n') for page in
    text_pages]))
3 text_lines = [line for line in text_lines if len(line) > 0]
4 if len(text_lines) == 0:
5     return []
6 model = fastbm25(text_lines)
7 result_keywords = model.top_k_sentence(keywords, k=3)
8 result_question = model.top_k_sentence(anoy_question, k=3)
9 top_match_indexes = [t[1] for t in result_question + result_keywords]
10 block_line_indexes = merge_idx(top_match_indexes, len(text_lines), 0, 30)
11
12 text_blocks = ['\n'.join([text_lines[idx] for idx in line_indexes]) for
    line_indexes in block_line_indexes]
13 # text_blocks = [filter_header_footer(text_block) for text_block in
    text_blocks]
14 text_blocks = [re.sub(' {3,}', '\t', text_block) for text_block in
    text_blocks]
15
16 text_blocks = [(t, SequenceMatcher(None, any_question, t,
    autojunk=False).find_longest_match().size) for t in text_blocks]
17 for text_block, match_size in text_blocks:
18     match = SequenceMatcher(None, any_question, text_block,
    autojunk=False).find_longest_match()
19     print(any_question[match.a: match.a + match.size])
20 max_match_size = max([t[1] for t in text_blocks])
21 text_blocks = [t[0] for t in text_blocks if t[1] == max_match_size]
22
23 if sum([len(t) for t in text_blocks]) > 2000:
24     max_avg_len = int(2000 / len(text_blocks))
25     text_blocks = [t[:max_avg_len] for t in text_blocks]
26
27 text_blocks = [rewrite_text_block(t) for t in text_blocks]
28 text_blocks = ['```\n{}\n```'.format(t) for t in text_blocks]
```

3. prompt_util.prompt_question_tp31组合召回的文本形成prompt传给模型进行回答