



KIST-CSRC / MaTableGPT



<> Code Issues Pull requests Actions Projects Security Insights

MaTableGPT / GPT\_models / follow\_up\_q.py



**gyeonghoonyi** fine tuning, few shot, zero shot, follow up questions

391089a · last year



598 lines (504 loc) · 35.5 KB

Code

## Blame

Raw



MaTableGPT/GPT\_models/follow\_up\_q.py at main · KIST-CSRC/MaTableGPT

```

answer strictly related to question 5. In the answer to question 4, remove any
parts corresponding to 'NA', 'na', 'unknown', or similar contents. Show the modified
31 JSON. Only display the [JSON representing "Question 4") based] on the input representation,
provide values of the {element} of the {catalyst_name} as a Python list. If there is a
32 unit, please provide strictly the value including the unit. The elements of a Python
list must be "Somewhat good" for Question 5. On the contrary, if the unit is "Very good"
33 Välgjord) then "]; value''' key of the sublayers of the {element} as a Python list. Only
output the Python list [{"Name": "Question 6": based}) only numerical values, provide a
34 list of elements that exist in the answer to Question 5 but are not present in the the
answer to Question 4: Note that "Question 4" is different from the elements in the
35 answer to Question 4: Note that "Question 4" is different from the elements in the
the answer to Question 4: Note that "Question 4" is different from the elements in the
36 answer to Question 4: Note that "Question 4" is different from the elements in the
37 answer to Question 4: Note that "Question 4" is different from the elements in the
# Questions related to electrolyte, reaction_type, substrate
self.questions_for_representation = {
    "1_str": ["representation_title", "Question 1. If there is any occurrence of
        'OER', 'HER', 'oxygen evolution reaction', 'hydrogen evolution reaction' or similar
        contents in the representation, does it correspond to the substrate? Please answer with either yes
        or no\n\n"], "4_str": ["representation_table_caption", "Question 4. Does the input
        representation contain information corresponding to substrate? Please answer with either yes or
        no\n\n"], "5_str": ["representation_title", "Question 5. Does the input representation
        contain information corresponding to electrolyte? Please answer with either yes or
        no\n\n"], "6_str": ["representation_table_caption", "Question 6. Does the input
        representation contain information corresponding to electrolyte? Please answer with
        either yes or no\n\n"]
}

# system prompt
self.system_prompt = {"role": "system", "content": "You need to modify the JSON
representing the table presenter.\n\n JSON template : {'catalyst_name' :
{'performance_name' : \{property template\}}}\n property template : {'electrolyte': '',
'reaction_type': '', 'value': '', 'current_density': '', 'overpotential': '',
'potential': '', 'substrate': '', 'versus': '', 'condition': ''}\n performance_list =
[overpotential, tafel_slope, Rct, stability, CdL, onset_potential, potential, TOF, ECSA,
water_splitting_potential, mass_activity, exchange_current_density, Rs,
51 specify type (self for overpotential, etc) loading,
52 apparent_activation_energy]\n In the JSON template, 'catalyst_name' and
53 'performance_name' as found in the repeated table header are present in the template.
representation."}

Parameters:
representation (str): The input representation string containing HTML content.
json (str): The JSON string to be included in the input.
want_type (str): The type of content required ('both', 'representation', 'json',
'representation_title', 'representation_table_caption').
Returns:
str: The formatted input string.
"""

```

```
63
64     # Split the representation string at each occurrence of '</table>'
65     splitted_strings = representation.split('</table>')
66     # Determine the format based on the requested type
67     if want_type == 'both':
68         format_for_input = "<input representation>\n" + str(representation) +
69     "n<input type='text' value='"
70         format_for_input = "<input representation>\n" + str(representation) + "\n\n"
71     elif want_type == 'json':
72         format_for_input = "<input json>\n" + str(json) + "\n\n"
73     elif want_type == 'representation_title':
74         title_string = splitted_strings[0] + '</table>'
75         format_for_input = "<input representation>\n" + str(title_string) + "\n\n"
76     elif want_type == 'representation_table_caption':
77         table_caption_string = splitted_strings[1]
78         format_for_input = "<input representation>\n" + str(table_caption_string) +
79     "\n\n"
80     return format_for_input
81
82     def load_file(self, file_type, file_path, file_name):
83     """
84         Loads the content of a file based on the specified type and path.
85
86         Parameters:
87             file_type (str): The type of the file (json, html, txt, etc.)
88             file_path (str): The path to the file
89             file_name (str): The name of the file
90
91         Returns:
92             output: The content of the file. The type of the content depends on the
93             file_type.
94         Raises:
95             ValueError: If the file type is unsupported.
96         """
97         with open(file_path + file_name, 'r', encoding='utf-8-sig') as f:
98             # Load JSON files
99             if file_type == 'json':
100                 output = json.load(f)
101             # Read text or HTML files
102             elif file_type in ['html', 'txt']:
103                 output = f.read()
104             # Raise an error for unsupported file formats
105             else:
106                 raise ValueError("Unsupported file format.")
107
108             return output
109
110     def formatting_type(self, key, answer):
111     """
112         Formats the answer based on the specified key type.
113
114         Parameters:
115             key (str): The key type for formatting (e.g., 'text', 'json', 'table').
```

```

115     key (str): The key indicating the type or format required (e.g., '1_list',
116     '2_dict'). answer (str): The answer to be formatted.
117
118     Returns:
119     The formatted answer, which can be a list or a dictionary based on the key type.
120     """
121
122     # Determine the desired format type from the key
123     want_type = key.split('_')[1]
124
125     # Handle formatting if the desired type is a list
126     if want_type == "list":
127         if answer[0] == "'" and answer[-1] == "'":
128             answer = answer.strip("'") # Remove surrounding quotes if present
129             answer = eval(answer) # Evaluate the string as a Python expression (e.g.,
130             convert to list)
131     # Handle formatting if the desired type is a dictionary
132     elif want_type == "dict":
133         if "```" in answer:
134             answer = answer.replace("```json", "").replace("```", "") # Remove
135             markdown code and work=fromat.loads(answer) # Parse the string as JSON
136
137     return answer
138
139     def check_type(self, key, answer):
140         """
141
142         Checks if the type of the answer matches the expected type based on the key.
143
144         Parameters:
145         key (str): The key indicating the expected type (e.g., '1_list', '2_dict').
146         answer (any): The answer whose type needs to be checked.
147
148         Returns:
149         tuple: A tuple containing the expected type (str) and a boolean indicating
150             whether the types match.
151         # Extract the question number and expected type from the key
152         question_number = key.split('_')[0]
153         want_type = key.split('_')[1]
154
155         # Determine the actual type of the answer
156         answer_type = type(answer).__name__
157
158         # Check if the expected type matches the actual type
159         type_bool = want_type == answer_type
160
161     def prompt(self, Q):
162         """
163
164         Sends a list of messages to the OpenAI ChatCompletion API and returns the
165         response.
166
167         Parameters:
168         Q (list): A list of message dictionaries to be sent to the API.

```

```
168         Returns:  
169             tuple: A tuple containing the original list of messages (Q) and the response  
170             content."  
171     while True:  
172         try:  
173             # Send a request to the OpenAI ChatCompletion API  
174             response = openai.ChatCompletion.create(  
175                 model="gpt-4-0125-preview",  
176                 messages=Q,  
177                 temperature=0,  
178                 frequency_penalty=0,  
179                 presence_penalty=0  
180             )  
181             break # Exit the loop if the request is successful  
182         except Exception as e:  
183             # Print the error message and retry  
184             print("An error occurred:", e)  
185  
186         # Return the original messages and the content of the response  
187         return Q, response['choices'][0]['message']['content']  
188  
189     def run(self, input_type, split_mode):  
190         """  
191             Executes the main process for modifying JSON files based on given questions.  
192  
193             Parameters:  
194             input_type (str): The type of input files (e.g., 'json', 'txt').  
195             split_mode (str): The mode for handling splits (e.g., 'split', 'no-split').  
196  
197             Returns:  
198             None  
199         """  
200         json_name_list = os.listdir(self.json_path)  
201         error_file_name = []  
202         transpose_bool = False  
203  
204         for json_name in json_name_list:  
205             # try:  
206                 txt_name = json_name.split('.')[0]  
207                 # Load input files  
208                 input_json = self.load_file('json', self.json_path, json_name)  
209                 input_representation = self.load_file(input_type,  
210                 self.representation_path, txt_name + '.txt')  
211                 # Initialize messages with system prompt  
212                 messages = [self.system_prompt]  
213  
214                 # Catalyst check  
215                 catalyst_result = {}  
216                 question_list = []  
217                 answer_list = []  
218                 question_token_list = []  
219                 answer_token_list = []
```

```

220     final_result = []
221     final_json = {"catalysts": []}
222
223     # Iterate through catalyst questions
224     for key, value in self.questions_for_catalyst.items():
225         if value[0] != "None":
226             messages.append({"role": "user", "content": ""})
227             self.input_prompt(messages.append(f"\nRole: {value[0]}"))
228             messages, answer = self.prompt(messages)
229             question_token_list.append(messages)
230             answer_token_list.append(answer)
231
232         try:
233             mod_answer = self.formatting_type(key, answer)
234         except Exception as e:
235             try:
236                 messages, answer = self.prompt(messages)
237                 question_token_list.append(messages)
238                 answer_token_list.append(answer)
239                 question_list.append("Throwing the same message once again")
240                 answer_list.append("Throwing the same message once again")
241             except Exception as e:
242                 error_file_name.append(json_name)
243                 print(json_name)
244                 break
245
246             catalyst_result[key[0]] = mod_answer
247             question_list.append(value[1])
248             answer_list.append(mod_answer)
249
250             print('-----')
251             print(self.questions_for_catalyst[key])
252             print(mod_answer)
253             print(catalyst_result)
254             print('-----')
255
256             messages.append({"role": "assistant", "content": answer})
257             if key[0] == '2' and catalyst_result["1"] == catalyst_result["2"]:
258                 print("GO TO THE PERFORMANCE STAGE !!")
259                 break
260
261             # Performance check
262             performance_result = {}
263             mod_catalyst_list = answer_list[-1]
264
265             if len(mod_catalyst_list) > 1:
266                 transpose_bool = True
267
268             for i in range(len(mod_catalyst_list)):
269                 messages = [self.system_prompt]
270                 for key, value in self.questions_for_performance.items():
271                     if value[0] != "None":
272                         messages.append({"role": "user", "content": ""})
273                         self.input_prompt(messages.append(f"\nRole: {value[0]}"))
274                         messages, answer = self.prompt(messages)
275                         question_token_list.append(messages)
276                         answer_token_list.append(answer)
277
278             catalyst_result[key[0]] = mod_answer
279             question_list.append(value[1])
280             answer_list.append(mod_answer)
281
282             print('-----')
283             print(self.questions_for_catalyst[key])
284             print(mod_answer)
285             print(catalyst_result)
286             print('-----')

```

```

273     self.input_prompt(input_representation=mod_catalyst_list[i], role = '""'+
274     mod_catalyst_list[i] +'Messages.append({\"role\": \"user\", \"content\": '+
275     messages, answer = self.prompt(messages)
276     question_token_list.append(messages)
277     answer_token_list.append(answer)
278
279     try:
280         mod_answer = self.formatting_type(key, answer)
281     except:
282         try:
283             messages, answer = self.prompt(messages)
284             question_token_list.append(messages)
285             answer_token_list.append(answer)
286         except:
287             error_file_name.append(json_name)
288
289         break
290
291     performance_result[f'{mod_catalyst_list[i]}_{key[0]}'] =
292     mod_answer
293     question_list.append(question)
294     answer_list.append(mod_answer)
295
296     messages.append({"role": "assistant", "content": answer})
297     if key[0] == '3' and
298     performance_result[f'{mod_catalyst_list[i]}_PROPERTY:STAGE !!']
299         break
300
301         if isinstance(answer_list[-1], list):
302             mod_performance_list =
303             performance_result[f'{mod_catalyst_list[i]}_1']
304             while not isinstance(answer_list[-1], list):
305                 messages.pop()
306                 messages, answer = self.prompt(messages)
307                 question_token_list.append(messages)
308                 answer_token_list.append(answer)
309                 mod_answer = self.formatting_type("1_list", answer)
310                 question_list.append(question)
311                 answer_list.append(mod_answer)
312                 mod_performance_list = answer
313
314             # Property check
315             property_result = {}
316             skip_questions = False
317             print("#####")
318             print(mod_performance_list)
319
320             if mod_performance_list == []:
321                 mod_answer = {str(mod_catalyst_list[i]): {}}
322             else:
323                 for j in range(len(mod_performance_list)):
324                     messages = []
325                     messages.append(self.system_prompt)

```

```

325         print("@@@@@@@@@@@@")
326         print(mod_performance_list[j])
327         # system prompt 넣어주고, input 표현자, json 넣어주는 코드
328         # messages = []
329         # messages.append(self.system_prompt)
330
331         skip_questions = False
332         answer3_no = False
333         for key, value in self.questions_for_property.items():
334             # 질문 2의 답변이 "no"일 때 질문 3번과 4번 skip하기 위한 조건문
335             if key[0] in ['3', '7'] and skip_questions:
336                 print("SKIP THE NEXT QUESTIONS")
337                 question_list.append("SKIP THE NEXT QUESTIONS")
338                 answer_list.append("SKIP THE NEXT QUESTIONS")
339                 if key[0] in ['3', '7']:
340                     skip_questions = False
341                 continue # 질문 3번과 4번 건너뛰기
342
343             if value[0] != "None":
344                 messages.append({"role": "user", "content": ""})
345                 self.input_prompt(input_representation, input_json, value[0]))
346                 question = value[1].format(catalyst_name = '""'+
347 mod_catalyst_list[i] + """', elementstage$append($performance$user$list{q}teht"?'question})'
348                 messages, answer = self.prompt(messages)
349                 question_token_list.append(messages)
350                 answer_token_list.append(answer)
351             try:
352                 mod_answer = self.formatting_type(key, answer)
353             except:
354                 try:
355                     messages, answer = self.prompt(messages)
356                     question_token_list.append(messages)
357                     answer_token_list.append(answer)
358                 except:
359                     print(json_name)
360                     error_file_name.append(json_name)
361                     break
362
363                 property_result[mod_catalyst_list[i] + '_' +
364 mod_performance_list[j] + '_'+ key[0]+json=list.append(question)
365                 answer_list.append(mod_answer)
366
367                 print('-----')
368                 print(question)
369                 print(mod_answer)
370                 print('-----')
371
372                 messages.append({"role": "assistant", "content": ""})
373                 answer})
374                 if key[0] == '2' and mod_answer.lower() == "no":
375                     question_list.append("Question 3. Based on the
answer to question 2, remove any parts corresponding to 'NA', 'na', 'unknown', or
...

```

```

3/6
377 similar content from the answer to question 1. Show the modified JSON. Only display the
378 answer_list.append(not(property_result[mod_catalyst_list[i]:+"user", "content":])
379 skip_questions = True
380 answer3_no = True
381
382 if key[0] == '6' and mod_answer == []:
383     if answer3_no:
384         question_list.append("Question 7. If the answer
385             to question 6 is an empty list, just provide the answer to question 1 as it is.")
386
387     answer_list.append(str(property_result[mod_catalyst_list[i]:+"user", "content":])
388     skip_questions = True
389 else:
390     question_list.append("Question 7. If the answer
391             to question 6 is an empty list, just provide the answer to question 3 as it is.")
392
393     answer_list.append(str(property_result[mod_catalyst_list[i]:+"user", "content":])
394     skip_questions = True
395
396     if isinstance(mod_answer, dict):
397         input_json = mod_answer
398     else:
399         count = 0
400         while not isinstance(mod_answer, dict) and count < 3:
401             messages.pop()
402             messages.append({"role": "user", "content"::
403 self.input_prompt(input_representation, answer['self'].prompt(messages)
404             question_token_list.append(messages)
405             answer_token_list.append(answer)
406             try:
407                 mod_answer = self.formatting_type(key, answer)
408             except:
409                 try:
410                     messages, answer = self.prompt(messages)
411                     question_token_list.append(messages)
412                     answer_token_list.append(answer)
413                 except:
414                     print('$$@#@#%#@##@#@#')
415                     print(json_name)
416                     print('$$@#@#%#@##@#@#')
417                     error_file_name.append(json_name)
418                     break
419
420             count += 1

```

```
421
422                         property_result[mod_catalyst_list[i] + '_' +
423     mod_performance_list[j] + '_' + key[0]]questions.append(question)
424                         answer_list.append(mod_answer)
425
426                         print('-----')
427                         print(question)
428                         print(mod_answer)
429                         print('-----')
430
431                         messages.append({"role": "assistant", "content": 
432     answer})
433             if len(mod_catalyst_list) == 1 and split_mode == 'split':
434                 final_result.append(mod_answer)
435             else:
436                 input_json = self.load_file(input_type, self.json_path,
437                     json_name)
438
439             if transpose_bool:
440                 final_result.append(final_json)
441
442             # Final JSON after property modifications
443             if final_result[0] == []:
444                 input_json = self.load_file('json', self.json_path, json_name)
445                 new_json = input_json
446             else:
447                 new_json = final_result[0]
448
449             # Handle representation questions
450             remove_elements = []
451             representation_result = {}
452             for key, value in self.questions_for_representation.items():
453                 messages = []
454                 if value[0] == "None":
455                     pass
456                 else:
457                     messages.append({"role": "user", "content": 
458             self.input_prompt(input_representation, new_json, value[0])})
459                     question = value[1]
460                     messages.append({"role": "user", "content": question})
461                     messages, answer = self.prompt(messages)
462                     question_token_list.append(messages)
463                     answer_token_list.append(answer)
464
465                     mod_answer = self.formatting_type(key, answer)
466
467                     question_list.append(question)
468                     answer_list.append(mod_answer)
469
470                     print('-----')
471                     print(question)
472                     print(mod_answer)
```

```

MaTableGPT/GPT_models/follow_up_q.py at main · KIST-CSRC/MaTableGPT
    473     print('-----')
    474
    475     messages.append({"role": "assistant", "content": answer})
    476     representation_result[key[0]] = mod_answer.replace(".", "").lower()
    477
    478     if representation_result['1'] == 'no' and representation_result['2'] ==
    479         'no':
    480             remove_elements.append('reaction_type')
    481     if representation_result['3'] == 'no' and representation_result['4'] ==
    482         'no':
    483             remove_elements.append('substrate')
    484     if representation_result['5'] == 'no' and representation_result['6'] ==
    485         'no':
    486             remove_elements.append('electrolyte')
    487
    488             remove_elements = list(set(remove_elements))
    489
    490
    491     if remove_elements != []:
    492         for delete_element in remove_elements:
    493             messages = []
    494             messages.append({"role": "user", "content": self.input_prompt(input_questions[removeawayelement] with the key name
    495             {delete_element} from the input JSON and display it in only JSON format. Other
    496             explanation is not allowed. message.append({JSONResultOnlydisplay:theJSON} like
    497             string not `` `json)` .format(delete_element= self.delete_element)`` ``"))
    498             question_token_list.append(messages)
    499             answer_token_list.append(answer)
    500             try:
    501                 mod_answer = self.formatting_type('1_dict', answer)
    502             except:
    503                 try:
    504                     messages, answer = self.prompt(messages)
    505                     question_token_list.append(messages)
    506                     answer_token_list.append(answer)
    507                 except:
    508                     print('#$@#@%#@##@#%#@#%#@#')
    509                     print(json_name)
    510                     print('#$@#@%#@##@#%#@#%#@#')
    511                     error_file_name.append(json_name)
    512                     break
    513
    514             question_list.append(question)
    515             answer_list.append(mod_answer)
    516             print('-----')
    517             print(question)
    518             print('answer')
    519             print(answer)
    520             print('mod answer')
    521             print(mod_answer)
    522             print('-----')
    523
    524     if not isinstance(mod_answer, dict):
    525         count = 0
    526         while not isinstance(mod_answer, dict) and count < 3:
    527             # 원래 했던 질문 제거하고 다시
    528             mod_answer = self.prompt(messages)

```

```

MaTableGPT/GPT_models/follow_up_q.py at main · KIST-CSRC/MaTableGPT
  messages, answer = self.prompt(messages)
    question_token_list.append(messages)
    answer_token_list.append(answer)
    try:
        mod_answer = self.formatting_type('1_dict', answer)
    except:
        try:
            messages, answer = self.prompt(messages)
            question_token_list.append(messages)
            answer_token_list.append(answer)
        except:
            print('#$@#@#%@##@#@#%@#')
            print(json_name)
            print('#$@#@#%@##@#@#%@#')
            error_file_name.append(json_name)
            break
    count += 1
    question_list.append(question)
    answer_list.append(mod_answer)

    print('-----')
    print(question)
    print(mod_answer)
    print('-----')

if not isinstance(mod_answer, dict):
    new_json = new_json
else:
    new_json = mod_answer
else:
    mod_answer = new_json

# Ensure the necessary directories exist
os.makedirs(os.path.join(self.save_path, 'log'), exist_ok=True)
os.makedirs(os.path.join(self.save_path, 'token'), exist_ok=True)
os.makedirs(os.path.join(self.save_path, 'json'), exist_ok=True)

# Save the modified JSON and log
if json_name not in error_file_name:
    log_path = self.save_path + 'log/' + txt_name
    df = pd.DataFrame({'Question': question_list, 'GPT answer':
answer_list})
    df.to_csv(log_path+'.csv', index=False)

    token_path = self.save_path + 'token/' + txt_name
    token_df = pd.DataFrame({'Question': question_token_list, 'GPT
answer': answer_token_list})
    token_df.to_csv(token_path+'.csv', index=False)

    new_json_path = self.save_path + 'json/' + json_name
    if mod_answer == []:
        input_json = self.load_file('json', self.json_path, json_name)
        with open(new_json_path, "w") as json_file:
            json.dump(input_json, json_file, indent=4)

```

```
577     else:
578         if isinstance(mod_answer, list):
579             with open(new_json_path, "w") as json_file:
580                 json.dump(new_json, json_file, indent=4)
581
582         elif isinstance(mod_answer, str):
583             with open(new_json_path, "w") as json_file:
584                 json.dump(new_json, json_file, indent=4)
585         else:
586             with open(new_json_path, "w") as json_file:
587                 json.dump(mod_answer, json_file, indent=4)
588
589
590     if __name__ == "__main__":
591         representation_path = 'table representer path'
592         json_path = 'gpt prediction'
593         save_path = 'save path'
594
595         assistant = FollowQ(json_path, representation_path, save_path)
596         assistant.run('txt', 'split')
```