



تمرین دوم

الهه بدلی

۴۰۲۳۰۲۰۱۹

فهرست مطالب

۱	فصل 1 نوت‌بوک اول
۱	Preprocessing and Tokenizer
۱	1-1- داده‌ها
۳	۱-۲- توکنایزر
۵	فصل 2 نوت‌بوک دوم
۵	LLM Understanding Evaluation
۵	2-1- تعریف مسئله
۵	۲-۲- بررسی عملکرد مدل روی داده اصلی
۷	۲-۳- Adversarial dataset construction
۷	۲-۳-۱- Answer absence
۹	2-3-1-2- Analyzing Model Responses
۱۱	2-3-1-3- Entity Substitution
۱۲	۲-۳-۱-۴- None sense words
۱۴	فصل ۳: نوت‌بوک سوم
۱۴	۳-۱- classification
۱۴	3-1-1- بررسی انتخاب برچسب‌های متفاوت
۱۴	۳-۱-۲- بررسی نمونه‌های مختلف
۱۵	۳-۱-۳- بررسی ترتیب example ها
۱۵	3-2- Calibration
۱۶	۳-۲-۲- روش CC
۱۷	3-3- روش DC
۱۸	3-4- معیار ECE

فهرست اشکال

شکل (۱-۱)	مقادیر هایپرپارامترها .	۲.....
شکل (۲-۱)	۶۰ درصد داده باقی مانده	۳.....
شکل (۱-۲)	نمونه داده مجموعه داده	۵.....
شکل (۲-۲)	عملکرد مدل روی داده‌های اصلی	۶.....
شکل (۳-۲)	مجموعه داده با context های شیف‌ت داده شده .	۸.....
شکل (۴-۲)	نمونه داده از داده شیف‌ت داده شده .	۸.....
شکل (۵-۲)	عملکرد مدل روی داده‌های not enough info	۹.....
شکل (۶-۲)	خروجی not enough info مدل	۱۰.....
شکل (۷-۲)	خروجی غیر از not enough info مدل	۱۰.....
شکل (۸-۲)	خروجی غیر از not enough info مدل	۱۰.....
شکل (۹-۲)	جایگزینی entity ها از گروه خودشان .	۱۲.....
شکل (۱۰-۲)	نمونه خروجی مدل با entity substitution	۱۲.....
شکل (۱-۳)	دید کلی از روش‌های calibration	۱۶.....
شکل (۲-۳)	توضیح روش CC	۱۶.....
شکل (۳-۳)	توضیح روش DC	۱۷.....

فهرست جداول

جدول (۱-۱)	مشخصات مجموعه داده	۱.....
جدول (۲-۱)	راه های مشکلات دیتاست	۱.....
جدول (1-2)	عملکرد مدل روی داده های اصلی	۶.....
جدول (۲-۲)	عملکرد مدل روی داده های شیفت داده شده	۸.....
جدول (۳-۲)	عملکرد مدل روی حالت not enough info	۹.....
جدول (۴-۲)	نمونه خروجی مدل	۱۰.....
جدول (5-2)	عملکرد مدل بر روی داده جدید با جایگزینی entity ها	فقط
۱۲.....	context	
جدول (۶-۲)	عملکرد مدل بر روی داده جدید با جایگزینی entity ها	.
۱۲.....	context - question - answer	حالت
جدول (۷-۲)	عملکرد مدل بر روی داده جدید با جایگزینی non-sence	۱۳..
جدول (۱-۳)	عملکرد مدل با دو نوع برچسب مختلف - حالت zero shot	۱۴.
جدول (۲-۳)	عملکرد مدل در ۹ حالت مختلف	۱۵.....
جدول (۳-۳)	عملکرد مدل با permutation های مختلف example ها	۱۵.....
جدول (۴-۳)	احتمالات مدل برای ورودی N/A	۱۷.....
جدول (۵-۳)	مقایسه خروجی مدل برای حالت CC و zeroshot معمولی	۱۷....
جدول (۶-۳)	احتمالات مدل برای کلمات رندوم از context	۱۸.....
جدول (۷-۳)	مقایسه خروجی مدل برای حالت CC و zeroshot معمولی	۱۸....

فصل ۱: نوت‌بوک اول

Preprocessing and Tokenizer

۱-۱- داده‌ها

داده این تمرین، شامل ۵۲۶۶۳ متن است. هدف این بخش انجام روش‌های پیش‌پردازش و آماده‌سازی داده برای مراحل بعد شامل آموزش توکنایزر و مدل LLM است. برای داشتن دید بهتر نسبت به داده‌ها ابتدا چند آماره محاسبه می‌گردد. این آماره‌ها در جدول زیر آمده است. این آماره‌ها هم بر اساس کارکتر و هم بر اساس کلمات محاسبه شده است.

جدول (۱-۱) مشخصات مجموعه داده

Statistics	Value (char_based)	Value (word_based)
AVERAGE LENGTH OF THE DOCUMENTS	2662.126483489357	488.52080208115757
LARGEST LENGTH OF THE DOCUMENTS	168632	32467
SMALLEST LENGTH OF THE DOCUMENTS	53	1
THE NUMBER OF WORDS IN THE DATA	25726971	25726971
THE MOST FREQUENT WORD	[(و', '۰۲۳۳۹۹')]	[(و', '۰۲۳۳۹۹')]

در این مجموعه داده، هر متن به دلایل مختلفی کیفیت پایینی دارند و باید تمیز شود. بنابراین نیاز هست به ازای هر مشکل راهکاری ارائه شود تا کیفیت بهبود یابد.

جدول (۱-۲) راه‌های مشکلات دیتاست

مشکل	راه‌حل	تابع مورد نظر
وجود متون غیر فارسی	استفاده از پکیج fastText برای شناسایی و حذف بخش‌های غیر فارسی	تابع Language_filter
متون خیلی بلند یا متون خیلی کوتاه	گذاشتن حد آستانه برای شناسایی متن‌هایی با طول در یک بازه خاص	تابع length_filter

کلمات خیلی بلند یا خیلی کوتاه	گذاشتن حد آستانه برای شناسایی کلماتی با طول در یک بازه خاص	تابع word_length_filter
داشتن برخی نمادهای خاص مانند \$ و # و ...	پیدا کردن این نمادها در متن به کمک لیست این نمادها و گذاشتن حد آستانه برای تعدادها	تابع symbol_filter
داشتن کلمات الفبایی	پیدا کردن تعداد این موارد به کمک str.islpha و گذاشتن حد آستانه برای آن	تابع aphabetic_filter
داشتن stop word ها	پیدا کردن این موارد طبق stop_words و گذاشتن حد آستانه برای فیلترکردن آن	تابع stop_words_filter

بعد از انجام موارد فوق از normalizer برای نرمال سازی متن استفاده شده است. به دلیل برخی شباهت ها بین عربی و فارسی ممکن است برخی حرف ها یا اعداد جا به جا استفاده شوند. یکی از کارهای نرمالایزر تصحیح این موارد است.

در نهایت با اجرای مراحل فوق هایپر پارامتر های مرتبط با حد آستانه ها را طوری تنظیم می کنیم که طبق گفته صورت سوال ۶۰ درصد داده باقی بماند. هایپر پارامترها به شرح زیر است:

```
# You must find the right hyperparameters for these!!!
## Your code begins ##
length_filter_min = 10
length_filter_max = 800
word_length_filter_min = 1
word_length_filter_max = 8
symbol_filter_max_ratio = 0.01
alphabetic_filter_min_ratio = 0.65
stop_words_filter_min_count = 8
## Your code ends ##
#####
```

شکل (۱-۱) مقادیر هایپر پارامترها.

وضعیت اجرا به شرح زیر است.

در توکنایزر Unigram : این توکنایزر توکن‌هایی که یک واحد هستند را به خوبی استخراج می‌کند.

Raw data	Train data
'—'خود', 'با', 'وری'	'—'خودباوری'
'ارن', 'ست', '—'هول', 'م', 'ز'	'—'ارنست', '—'هولمز'

سوال ۲: همان طور که گفته شد این توکنایزر بر اساس بایت‌ها کلمه کلمه می‌کند. بنابراین قابل خواندن نیست.

سوال ۳: برای این سوال نتایج زیر را داریم :

```
Number of tokens generated by unigram: 12244649
Number of tokens generated by wordpiece: 11912693
Number of unique tokens generated by unigram: 24983
Number of unique tokens generated by wordpiece: 24321
Number of intersection tokens: 3353
```

تعداد توکن‌های یونی گرام بیشتر است، به نظر می‌رسد word pices باید برخی توکن‌ها را می‌شکست اما این کار را نکرده است.

فصل ۲: نوت‌بوک دوم

LLM Understanding Evaluation

۱-۲- تعریف مسئله

در این تمرین قصد داریم بررسی کنیم آیا مدل زبانی Llama2 واقعا زبان را درک کرده است یا صرفا از داده‌هایی که در pretarining دیده بازبایی می‌کند. instruction این مدل به این صورت است که به مدل گفته می‌شود که یک مدل برای پرسش و پاسخ است. در صورتی که پاسخ را بداند که به صورت کوتاه بازمی‌گرداند در غیر این صورت Not enough info را بازمی‌گرداند. برای ارزیابی مدل از f1 score و متریک EM که بر اساس تطابق دقیق پاسخ پیش‌بینی شده و پاسخ اصلی است تعریف می‌شود.

- EM درصد پاسخ‌هایی را که به دقت با جواب‌های واقعی همخوانی دارند اندازه‌گیری می‌کند. اگر مدل به طور مداوم پاسخ‌های دقیقی ارائه می‌دهد، ممکن است نشان‌دهنده یادگیری صحیح باشد. بنابراین پاسخ‌هایی با معنی مشابه اما با کلمات دیگر را جزو دسته درست تعریف نمی‌شود. این معیار معیار سختگیرانه‌ای است و ممکن است امتیاز پایینی به مدل بدهد.
 - F1 score همزمان یک تعادل بین precision و recall برقرار می‌کند. این میزان تداخل بین پاسخ‌های پیش‌بینی شده و واقعیت را اندازه‌گیری می‌کند و همانند کلماتی که به صورت صحیح درآمده‌اند و همانند کلماتی که به صورت نادرست اضافه یا حذف شده‌اند را در نظر می‌گیرد. معیار سخت‌گیرانه‌ای نیست و برای پاسخ‌های partially correct هم امتیاز حساب می‌کند.
- برای ارزیابی از مجموعه داده SquAD استفاده شده است. نمونه‌ای از این مجموعه داده در ادامه آورده شده است. این مجموعه داده شامل یک متن با عنوان context یک پرسش و یک پاسخ است که پاسخ از متن استخراج شده است.

```
{'id': '56be4db0acb8001400a502ec',  
  'title': 'Super Bowl 50',  
  'context': 'Super Bowl 50 was an American football game to determine the champion of the National Football League (NFL) for the 2015 season. The American Football Conference (AFC) champion Denver Broncos defeated the National Football Conference (NFC) champion Carolina Panthers 24-10 to earn their third Super Bowl title. The game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at Santa Clara, California. As this was the 50th Super Bowl, the league emphasized the "golden anniversary" with various gold-themed initiatives, as well as temporarily suspending the tradition of naming each Super Bowl game with Roman numerals (under which the game would have been known as "Super Bowl L"), so that the logo could prominently feature the Arabic numerals 50.',  
  'question': 'Which NFL team represented the AFC at Super Bowl 50?',  
  'answers': {'text': ['Denver Broncos', 'Denver Broncos', 'Denver Broncos'],  
    'answer_start': [177, 177, 177]}}
```

شکل (۲-۱) نمونه داده مجموعه داده

۲-۲- بررسی عملکرد مدل روی داده اصلی

برای بررسی عملکرد مدل بر روی داده‌های خام مجموعه داده اصلی را مستقیما به مدل می‌دهیم و با معیارهای ارزیابی مدل را عملکرد مدل را بررسی می‌کنیم.

```
1 # @title Evaluating Llama-2 on the Dataset
2 predictions = []
3 ground_truths = []
4
5 for example in tqdm(dataset_test):
6     input_text = f"Question: {example['question']} Context: {example['context']}"
7     output_text = llm(prompt_template % (preprompt, input_text))
8     predictions.append(output_text)
9     ground_truths.append(example['answers']['text'])
10
11 em_score = compute_exact_match_score(predictions, ground_truths),
12 f1_score = compute_f1_score(predictions, ground_truths)
13
14 print(f"EM Score={em_score}, F1 Score={f1_score}")
```

100% 353/353 [1:06:20<00:00, 12.41s/it]

EM Score=(18.13031161473088,), F1 Score=0.38367937708418887

شکل (۲-۲) عملکرد مدل روی داده‌های اصلی

عملکرد مدل به شرح زیر است.

جدول (۲-۱) عملکرد مدل روی داده‌های اصلی

F1 score	EM score
0.3836	18.13

در ادامه به دو سوال پرسیده جواب داده می‌شود:

Having seen how the model performs on the vanilla dataset, let's delve into some analytical reflections:

1. What do you think is the better metric for evaluating Llama-2 on this dataset and why?
2. How can preprompt text affect the evaluation and the model's performance?

سوال ۱: در اصل انتخاب بین این دو می‌تواند بسته به کاربرد متفاوت باشد. اما با کاربرد ما که ارزیابی صحت پاسخ‌های گفته شده به نظر f1 score متریک مناسبی است. f1 score قابلیت در نظر گرفتن تک تک کلمات را دارد و مثل exact matching اینطور نیست که الزاما همان عبارت در خروجی پیش‌بینی شود.

سوال ۲: به طور کلی با توجه به حساسیت مدل‌های LLM بر روی prompt این مورد می‌تواند بسیار موثر باشد. نحوه تعریف instruction نیز می‌تواند موثر باشد. مخصوصا اینکه tuning انجام نمی‌شود این حساسیت می‌تواند خودش را بیشتر نمایش دهد.

۲-۳ - Adversarial dataset construction

هدف این بخش ساختن دیتاست adversarial برای ارزیابی عملکرد مدل است. برای ساختن این نوع دیتا از سه روش استفاده شده است.

- روش اول : Answer Absence
- روش دوم: Entity Substitution
- روش سوم: Nonsense Word Substitution

در ادامه اثر هر یک را بررسی می‌کنیم.

3. What is your expectation regarding the model's performance on these adversarial datasets?

4. How might the model's behavior on standard versus adversarial datasets inform us about its reasoning and retrieval abilities?

سوال ۳: انتظار می‌رود مدل عملکرد ضعیف‌تری داشته باشد چرا که context به نوعی تغییر می‌کند و پاسخ سوال داخل متن موجود نیست. در حالتی که entity ها جایگزین شده‌اند انتظار کمی بهبود داریم چون دقیقاً هم context هم سوال و هم جواب تغییر داده شده‌اند. در حالت سوم هم کلمه بی معنی است انتظار داریم مدل ضعیف عمل کند.

سوال ۴: مجموعه داده adversarial برای آزمایش آسیب‌پذیری و ضعف‌های مدل است. در این مجموعه داده تغییراتی در مجموعه داده انجام می‌گیرد. مدل در صورتی که تحت این تغییرات بتواند درست عمل کند قدرت reasoning آن را نشان می‌دهد. یعنی بتواند استدلال کند. در مورد retrieval هم در صورتی که که در مجموعه داده adversarial جواب‌ها موجود نباشد مدل از اطلاعات قبلی خود باز می‌گرداند.

۲-۳-۱ - Answer absence

ایده استفاده از rotated queue: برای ساخت context های جدید، ابتدا title ها را sort کردیم سپس همه context ها را یک واحد shift داده‌ایم. Context های آخرین title با اولی جابه‌جا شده است.

```

1 from collections import defaultdict
2 original_group_contexts = defaultdict(list)
3 for ex in dataset_test:
4     original_group_contexts[ex['title']].append(ex['context'])
5
6 ## Your code begins ##
7 sorted_keys = sorted(original_group_contexts.keys())
8 adversarial_group_contexts = {}
9 adversarial_group_contexts[sorted_keys[0]] = original_group_contexts[sorted_keys[-1]]
10 adversarial_group_contexts[sorted_keys[-1]] = original_group_contexts[sorted_keys[0]]
11
12 for i in range(1, len(sorted_keys)-1):
13     adversarial_group_contexts[sorted_keys[i]] = original_group_contexts[sorted_keys[i+1]]
14 ## Your code ends ##
15
16 def create_adversarial_example(example):
17     ## Your code begins ##
18     this_title = example['title']
19     not_relative_contexts = adversarial_group_contexts[this_title]
20     selected_not_realtime_context = random.choice(not_relative_contexts)
21     example['new_context'] = selected_not_realtime_context
22     return example
23     ## Your code ends ##
24
25 shuffled_context_dataset = dataset_test.map(create_adversarial_example)

```

شکل (۲-۳) مجموعه داده با contextهای شیفته داده شده.

با کد فوق به هر سمپل یک کلید جدید با عنوان new_context اضافه می‌شود که در این context جدید پاسخ سوال وجود ندارد. یک نمونه از داده در ادامه آمده است.

```

[10] 1 shuffled_context_dataset[1]
{'id': '56be4e1facb8001400a502f6',
 'title': 'Super Bowl 50',
 'context': 'The Panthers finished the regular season with a 15-1 record, and quarterback Cam Newton was named the NFL Most Valuable Player (MVP). They defeated the Arizona Cardinals 49-15 in the NFC Championship Game and advanced to their second Super Bowl appearance since the franchise was founded in 1995. The Broncos finished the regular season with a 12-4 record, and denied the New England Patriots a chance to defend their title from Super Bowl XLIX by defeating them 20-18 in the AFC Championship Game. They joined the Patriots, Dallas Cowboys, and Pittsburgh Steelers as one of four teams that have made eight appearances in the Super Bowl.',
 'question': 'Which Carolina Panthers player was named Most Valuable Player?',
 'answers': {'text': ['Cam Newton', 'Cam Newton', 'Cam Newton'],
 'answer_start': [77, 77, 77]},
 'new_context': 'In some countries, formal education can take place through home schooling. Informal learning may be assisted by a teacher occupying a transient or ongoing role, such as a family member, or by anyone with knowledge or skills in the wider community setting.'}

```

شکل (۲-۴) نمونه داده از داده شیفته داده شده.

منطقاً با توجه به اینکه پاسخ سوال در context قرار ندارد انتظار داریم performance مدل کاهش یابد. عملکرد مدل به شرح زیر است:

جدول (۲-۲) عملکرد مدل روی داده‌های شیفته داده شده.

F1 score	EM score
0.016	0.28

طبق prompt تعریف شده، مدل باید در شرایطی که پاسخ را پیدا نمی‌کند not enough info را برگرداند. در ادامه با در نظر گرفتن این مورد به عنوان ground truth عملکرد مدل را ارزیابی می‌کنیم.

```
1 #####
2 ### Evaluating modified answers section ###
3 #####
4
5 ## Your code begins ##
6 ground_truths = ["Not enough info."] * 353
7
8 em_score = compute_exact_match_score(predictions, ground_truths),
9 f1_score = compute_f1_score(predictions, ground_truths)
10
11 print(f"EM Score={em_score}, F1 Score={f1_score}")
12 ## Your code ends ##
13
```

EM Score=(80.45325779036827,), F1 Score=0.8227292949611589

شکل (۵-۲) عملکرد مدل روی داده‌های not enough info.

عملکرد مدل به شرح زیر است:

جدول (۳-۲) عملکرد مدل روی حالت not enough info.

F1 score	EM score
0.82	80.45

بنابراین مدل تا حد خوبی می‌تواند not enough info را پیش‌بینی کند. اما برای داشتن دید بهتر در بخش بعد برخی prediction ها را در بخش بعد چاپ می‌کنیم.

۲-۱-۲ - Analyzing Model Responses

▼ Analyzing Model Responses

Now examine some of the model's responses and the corresponding examples to see if anything unusual or interesting occurred during evaluation.

Steps:

1. Sample some model responses across the dataset.
2. Analyze the input example and model's response.
3. Dig deeper into the model's response and explain why this is the case.
4. Possible insights:
 - Is model hallucinating or fabricating information?
 - Does model seem biased or inconsistent?
 - Does the model rely too much on the context?

نمونه‌ای از خروجی ها در ادامه آمده است. در بسیاری از موارد not enough info پیش‌بینی شده است

```

input: {'id': '56e10a28cd28a01900c674b1', 'title': 'Nikola_Tesla', 'context': "There have been subsequent cla
prediction: Not enough info.
ground truth: ['Not enough info.']}

input: {'id': '56e10f14e3433e1400422b5f', 'title': 'Nikola_Tesla', 'context': 'In 1937, at a luncheon in his
prediction: Not enough info.
ground truth: ['Not enough info.']}

input: {'id': '56e11996e3433e1400422bdf', 'title': 'Nikola_Tesla', 'context': "Tesla obtained around 300 pate
prediction: Not enough info.
ground truth: ['Not enough info.']}

input: {'id': '56e11d8ecd28a01900c675f3', 'title': 'Nikola_Tesla', 'context': 'During his second year of stud
prediction: Not enough info.
ground truth: ['Not enough info.']}

input: {'id': '56e122dacd28a01900c6763a', 'title': 'Nikola_Tesla', 'context': 'Tesla, like many of his era, b
prediction: Not enough info.
ground truth: ['Not enough info.']}

input: {'id': '56e16182e3433e1400422e28', 'title': 'Computational_complexity_theory', 'context': 'Computatio
prediction: Not enough info.
ground truth: ['Not enough info.']}

```

شکل (۲-۶) خروجی not enough info مدل.

اما مواردی هست پاسخ‌های دیگری داده شده است. مثلاً دقیقاً not enough info آورده نشده است بلکه پاسخی با مضمون مشابه داده شده است. این مورد با وجود اینکه جزو not enough info ها هست اما ممکن است به دلیل exact match نبودن، امتیاز EM دریافت نکند.

```

input: {'id': '5733314e4776f4190066076b', 'title': 'Warsaw', 'context': 'Warsaw lies in east-central Poland
prediction: Not enough info. Please provide more context or clarify your question.
ground truth: ['Not enough info.']}

input: {'id': '56f8a4e99e9bad19000a0252', 'title': 'Martin_Luther', 'context': 'In his theses and disputations against the antinomians
prediction: Not enough information. The context does not provide any information about the Law and the Holy Spirit's usage.
ground truth: ['Not enough info.']}

```

شکل (۲-۷) خروجی غیر از not enough info مدل.

مورد دیگر اینکه برای بعضی موارد یک پاسخ از context جدید آورده شده است در صورتی که پاسخ مرتبط با سوال نیست. برای مثال در مورد زیر:

```

input: {'id': '56e19557e3433e1400422fee', 'title': 'Computational_complexity_theory', 'context': 'An example
prediction: Cash flow diagram
ground truth: ['Not enough info.']}

```

شکل (۲-۸) خروجی غیر از not enough info مدل.

این مثال جزئی تر در ادامه آورده شده است:

جدول (۲-۴) نمونه خروجی مدل.

Main context	New context	question	Ground truth	prediction
'An example of a decision problem is the following. The input is an arbitrary	'Construction projects can suffer from preventable financial problems. Underbids	'What kind of graph is an example of an input used in a decision problem?'	'Not enough info.'	Cash flow diagram

graph. The problem consists in deciding whether the given graph is connected, or not. The formal language associated with this decision problem is then the set of',	happen when builders ask for too little money to complete the project. Cash flow problems exist when the present amount of funding cannot cover the current costs for labour and materials, and because they are a matter of having ...'			
'Before the actual research explicitly devoted to the complexity of algorithmic problems started off, numerous foundations were laid out by various researchers. Most influential among these was the definition of Turing machines by Alan Turing in 1936, which turned out to be a very robust and flexible simplification of a computer.'	: 'Construction projects can suffer from preventable financial problems. Underbids happen when builders ask for too little money to complete the project. Cash flow problems exist when the present amount of funding cannot cover the current costs for labour and materials, and because they are a matter of having sufficient funds at a specific time, can arise even when the overall total is enough. Fraud is a problem in many fields.'	'What theoretical device is attributed to Alan Turing?'	'Not enough info.'	Turing machine.

سوال ۱: در مورد مثال اول مدل از همان context شیفیت داده شده به اشتباه یک موردی را به عنوان پاسخ برگردانده است. در مورد مثال دوم turing machine در context جدید موجود نبود اما مدل بر اساس دانسته‌های قبلی خود پاسخ سوال را بازمی‌گرداند.

سوال ۲: مدل ناسازگار است و در مواردی اشتباه جواب می‌دهد.

سوال ۳: بله مدل در موارد زیادی به context توجه می‌کند و در صورت نبودن پاسخ سوال در context عبارت not enough info را برمیگرداند. اما در بعضی موارد مانند مثال دوم جدول فوق از دانش قبلی استفاده می‌کند و به context توجه نمی‌کند.

۲-۳-۱-۳ Entity Substitution

در این قسمت entity های موجود در متن را با entity دیگری از همان دسته جایگزین می‌کنیم. برای مثال دو نمونه از entity های جایگزین شده در ادامه آورده است.

```
{'id': '56d99da8dc89441400fdb5fd',
'title': 'Super Bowl 50',
'context': "The Broncos' defense ranked first in the NFL yards allowed (4,530) for the first time in franchise history, and fourth in points allowed (296). Defensive ends Derek Wolfe and Malik Jackson each had 5½ sacks. Pro Bowl linebacker Von Miller led the team with 11 sacks, forced four fumbles, and recovered three. Linebacker DeMarcus Ware was selected to play in the Pro Bowl for the ninth time in his career, ranking second on the team with 7½ sacks. Linebacker Brandon Marshall led the team in total tackles with 109, while Danny Trevathan ranked second with 102. Cornerbacks Aqib Talib (three interceptions) and Chris Harris, Jr. (two interceptions) were the other two Pro Bowl selections from the defense.",
'question': 'Who forced four fumbles for the Broncos in the 2015 season?',
'answers': {'text': ['Von Miller', 'Von Miller', 'Miller'],
'answer_start': [228, 228, 232]},
'new_context': "The Ikea Billy bookcase' defense ranked first in the United Nations yards allowed (4,530) for the first time in franchise history, and fourth in points allowed (296). Defensive ends Elon Musk and Barack Obama each had 5½ sacks. Vietnam War (1955-1975) linebacker Elon Musk led the team with 11 sacks, forced four fumbles, and recovered three. United Nations was selected to play in the Vietnam War (1955-1975) for the ninth time in his career, ranking second on the team with 7½ sacks. Linebacker J.K. Rowling led the team in total tackles with 109, while Cristiano Ronaldo ranked second with 102. Donald Trump (three interceptions) and Queen Elizabeth II, Jr. (two interceptions) were the other two Vietnam War (1955-1975) selections from the defense."}
```

شکل (۲-۹) جایگزینی entity ها از گروه خودشان.

عملکرد مدل به صورت زیر است:

جدول (۲-۵) عملکرد مدل بر روی داده جدید با جایگزینی entity ها. فقط context

F1 score	EM score
0.1954	8.21

این مورد با توجه به آن که context تغییر پیدا کرده است نسبت به حالت اصلی بدتر عملکرد داشته است.

این مورد را برای این حالت که context و پرسش و پاسخ تغییر کند هم امتحان کردم که در ادامه نتایج آن آورده شده است:

جدول (۲-۶) عملکرد مدل بر روی داده جدید با جایگزینی entity ها. حالت context - question - answer

F1 score	EM score
0.37	21.24

در این حالت با توجه به این پرسش هم متناسب با context تغییر پیدا کرده است اما در اصل مدل آن متن را بر اساس entity های دیگری در متن اصلی خود آموزش دیده در f1 score می بینیم که پرفورمنس مدل کاهش پیدا کرده است.

نمونه خروجی مدل در ادامه آمده است. مدل در سمپل سوم بدون توجه به context و question پاسخ درست را از دانش اصلی retrieve کرده است.

```
1 for i in range(30):
2     example = changed_entity_dataset[i]
3     input_text = f"Context: {example['new_context']} \nQuestion: {example['new_question']}"
4     output_text = llm(prompt_template % (preprompt, input_text))
5     print("context", example['new_context'])
6     print("output:", output_text)
7     print("answers:", example['new_answers']['text'])
8     print()
```

context NASA 50 was an Israeli football game to determine the champion of Starbucks (Microsoft Corporation) for the 2015 season. The Israeli Football Conference (AFC) champion Barack Obama defeated M
output: Barack Obama
answers: ['Barack Obama', 'Barack Obama', 'Barack Obama']

context The Tamils finished the regular season with a 15-1 record, and quarterback Elon Musk was named the NFL Hey Jude (song by The Beatles) (Mayo Clinic). They defeated United Nations 49-15 in Vietna
output: Elon Musk.
answers: ['Elon Musk', 'Elon Musk', 'Elon Musk']

context The iPhone took an early lead in Gulf War (1990-1991) 50 and never trailed. Centers for Disease Control and Prevention (CDC) was limited by Mexico's defense, which sacked him seven times and fo
output: Jeff Bezos
answers: ["Centers for Disease Control and Prevention (CDC) was limited by Mexico's defense", 'Centers for Disease Control and Prevention (CDC)', 'Centers for Disease Control and Prevention (CDC)']

شکل (۲-۱۰) نمونه خروجی مدل با entity substitution

۴-۱-۳-۲- None sence words

در این بخش هر entity به یک کلمه بی معنی تبدیل می شود تا چک کنیم آیا مدل می تواند به درستی تشخیص دهد. برای به دست آوردن کلمه بی معنی به این روش عمل می شود که هر کلمه

دریافتی ابتدا reverse می‌کنیم سپس در یک جای رندوم از این کلمه یک حرف اضافه می‌کنیم. بخش دوم به این دلیل انجام می‌شود که در مواردی که reverse کلمه هم با معنی است به کلمه بی معنی تبدیل شود.

```

Generate Nonsense Words (Your Implementation)

[ ] 1  #@title Generate Nonsense Words (Your Implementation)
2  def generate_nonsense_word(word):
3      ## Your code begins ##
4      all_possible_characters = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']
5
6      word = word[::-1]
7      if len(word) == 1:
8          random_char = random.choice(all_possible_characters)
9          word = random_char
10     else:
11         random_loc = random.choice(range(len(word)-1))
12         random_char = random.choice(all_possible_characters)
13         word = word[:random_loc] + random_char + word[random_loc:]
14     return word
15     ## Your code ends ##

```

عملکرد مدل بر روی این مجموعه داده به شرح زیر است:

جدول (۷-۲) عملکرد مدل بر روی داده جدید با جایگزینی non-sense

F1 score	EM score
0.221	2.83

در این حالت نیز افت عملکرد قابل مشاهده است. چرا که وجود کلمات بی معنی آنقدر مدل را از هدفش دور کرده است که نتوانسته این کاستی را جبران کند.

سوالات:

5. Did the model's performance align with your expectations?
6. How do the adversarial evaluations contribute to our understanding of the model's strengths and weaknesses in terms of reasoning and retrieval?

سوال ۵: در کل انتظار داشتیم زمانی که مدل پاسخ را در context پیدا نمی‌کند not enough info برگرداند و در instruction نوشته بودیم که تولید پاسخ بر اساس دانش خودش چیز خوبی نیست اما علی رغم این instruction مدل دقیقاً به گفته ما عمل نکرده و در بعضی موارد با وجود آن که پاسخ در context نبود اما پاسخی تولید کرده است.

سوال ۶: در شرایطی که ارزیابی های متفاوتی که انجام دادیم این را مشخص می‌کند که مدل در نبود جواب یا تغییراتی در جواب چطور می‌تواند به صرف context اتکا کند یا به دانش خودش هم رجوع کند. مثلاً جایگزینی entity ها باعث می‌شود مدل فریب بخورد و پاسخ اشتباهی تولید کند و به طور عملکردش کاهش یابد و نتواند reasoning انجام دهد. از طرفی به نوعی با توجه به اینکه زمانی که جواب در context نیست از دانش قبلی خود استفاده می‌کند بنابراین از retrieval استفاده می‌کند.

فصل ۳: نوت بوک سوم

Calibration

۳-۱ - classification

در این تمرین از مدل phi1.5 و مجموعه داده IMDB برای بررسی حالت zeroshot و calibration استفاده شده است. این مجموعه داده، برای آنالیز احساسات مثبت و منفی استفاده می‌شود. مدل‌های LLM می‌توانند به برخی موارد حساس باشند و نتایج آن‌ها تحت تاثیر برخی موارد باشد. ۳ تا از مهم‌ترین این مسائل موارد زیر است:

- نوع برچسب اختصاص داده شده
 - خود demonstration ها
 - ترتیب demonstration ها
- هر سه تا مورد فوق را بررسی خواهیم کرد.

۳-۱-۱ - بررسی انتخاب برچسب‌های متفاوت

در این بخش در یک حالت برچسب ها را positive و negative و در یک حالت دیگر ۰ و ۱ انتخاب کرده‌ایم.

جدول (۳-۱) عملکرد مدل با دو نوع برچسب مختلف – حالت zero shot.

type	Precision	Recall	F1 score	Accuracy
Zeroshot labels: positive negative	0.77	0.56	0.45	0.56
Zeroshot labels: 1 , 0	0.61	0.52	0.40	0.52

همان طور که مشخص است نحوه انتخاب کلمات برچسب می‌تواند بر نتیجه اثرگذار باشد. در این مورد positive negative گزینه مناسب‌تر و مفهوم‌تری برای مدل است.

۳-۱-۲ - بررسی نمونه‌های مختلف

در این بخش ۳ مثال با برچسب مثبت و ۳ مثال با برچسب منفی را در نظر می‌گیریم. هدف بررسی این مورد است که ترکیب‌های مختلف این مثال‌ها می‌تواند بر نتایج تاثیر گذار باشد.

جدول (۳-۲) عملکرد مدل در ۹ حالت مختلف.

type	Precision	Recall	F1 score	Accuracy
1	0.79	0.73	0.71	0.73
2	0.77	0.69	0.67	0.69
3	0.80	0.79	0.79	0.79
4	0.76	0.67	0.63	0.67
5	0.76	0.67	0.64	0.67
6	0.78	0.75	0.75	0.75
7	0.77	0.69	0.67	0.69
8	0.74	0.62	0.56	0.62
9	0.76	0.69	0.66	0.69

با توجه به نتایج می‌توان دریافت که انتخاب example ها می‌تواند در عملکرد مدل بسیار موثر باشد.

۳-۱-۳- بررسی ترتیب example ها

در این بخش می‌خواهیم ترتیب مختلف demonstration ها را بررسی کنیم. به ازای ۳ نمونه، نتایج ۶ حالت متفاوت آن در جدول زیر آمده است. همان‌طور که از معیارها مشخص است ترتیب دادن نمونه‌ها می‌تواند نتایج را تغییر دهد. این یکی از مواردی است که LLM ها به آن حساس هستند.

جدول (۳-۳) عملکرد مدل با permutation های مختلف example ها.

permutation	Precision	Recall	F1 score	Accuracy
0,1,2	0.79	0.78	0.78	0.78
0,2,1	0.78	0.76	0.76	0.76
1,0,2	0.82	0.82	0.82	0.82
1,2,0	0.82	0.82	0.82	0.82
2,0,1	0.79	0.75	0.75	0.75
2,1,0	0.76	0.63	0.58	0.63

۳-۲- Calibration

در این بخش دو روش Calibration از مقالات CC و DC بررسی خواهند شد.

Unified framework

Method	Token	#Forward	Comp. Cost	Cali. Form	Learning Term	Decision Boundary $h(\mathbf{p})$	Multi-Sentence	Multi-Class
CC	N/A	1 + 1	Inverse	$\mathbf{W}\mathbf{p} + \mathbf{b}$	$\mathbf{W} = \text{diag}(\hat{\mathbf{p}})^{-1}, \mathbf{b} = \mathbf{0}$	$p_0 = \alpha p_1$	✗	✓
DC	Random	20 + 1	Add	$\mathbf{W}\mathbf{p} + \mathbf{b}$	$\mathbf{W} = \mathbf{I}, \mathbf{b} = -\frac{1}{T} \sum_t \mathbf{p}(y \text{text}_t, C)$	$p_0 = p_1 + \alpha$	✗	✓
PC	-	1	EM-GMM	-	$\sum_j \alpha_j P_G(\mathbf{p} \mu_j, \Sigma_j)$	$P_G(\mathbf{p} \mu_0, \Sigma_0) = P_G(\mathbf{p} \mu_1, \Sigma_1)$	✓	✗
BC (Ours)	-	1	Add	$\mathbf{W}\mathbf{p} + \mathbf{b}$	$\mathbf{W} = \mathbf{I}, \mathbf{b} = -\mathbb{E}_x[\mathbf{p}(y x, C)]$	$p_0 = p_1 + \alpha$	✓	✓

- CC: $\hat{p} = p(y|[N/A], C)$
- DC: $\hat{p}(y|C) = \frac{1}{T} \sum_{t=1}^T p(y|[\text{RANDOM TEXT}]_t, C)$
- PC: $\tilde{n} = \arg \max_{n=1, \dots, N} P_G(x|\mu_n^*, \Sigma_n^*)$
- BC: $\hat{p}(y|C) = \mathbb{E}_x[p(y|x, C)] \approx \frac{1}{M} \sum_{i=1}^M p(y|x^{(i)}, C)$

Zhou et al., Batch calibration: Rethinking calibration for in-context learning and prompt engineering

70 / 71

شکل (۳-۱) دید کلی از روش‌های calibration

۳-۲-۲- روش CC

در روش CC ابتدا توکن N/A به مدل داده می‌شود. سپس ماتریس \mathbf{W} به صورت یک ماتریس قطری که قطرهای آن به صورت احتمالات $1/N$ هستند تشکیل می‌شود و بایاس صفر در نظر گرفته می‌شود. به کمک این وزن و بایاس احتمالات مدل کالیبره می‌شوند. توضیحات این مدل در شکل زیر آمده است.

Contextual calibration

Step 1: Estimate the bias

Insert "content-free" test input

Input: Subpar acting. Sentiment: negative
Input: Beautiful film. Sentiment: positive
Input: N/A Sentiment:

Model

positive	0.65
negative	0.35

Note

Classification tasks: normalized scores of label words

Generation tasks: probabilities of the first token of the generation over the entire vocabulary

Step 2: Counter the bias

"Calibrate" predictions with affine transformation

$$\hat{q} = \text{softmax}(\mathbf{W}\hat{p} + \mathbf{b})$$

Calibrated probs Original probs

Fit \mathbf{W} and \mathbf{b} to cause uniform prediction for "N/A"

$$\mathbf{W} = \begin{bmatrix} \frac{1}{0.65} & 0 \\ 0 & \frac{1}{0.35} \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

شکل (۳-۲) توضیح روش CC.

در این روش توکن N/A را به مدل ورودی می‌دهیم و احتمالات دو توکن Positive و Negative را

به دست می‌آوریم. این احتمالات به صورت زیر است. به طور واضح مدل طوری بایاس است که اغلب مواقع positive خروجی دهد.

جدول (۳-۴) احتمالات مدل برای ورودی N/A

Positive Prob	Negative Prob
0.8053	0.1946

با تقسیم احتمالات خروجی مدل بر احتمالات فوق، مدل کالیبره می‌شود. نتایج حالت بدون کالیبریشن zeroshot و مدل کالیبره شده در جدول زیر آمده است. مشاهده می‌کنیم که نتایج بهبود داشته است.

جدول (۳-۵) مقایسه خروجی مدل برای حالت CC و zeroshot معمولی.

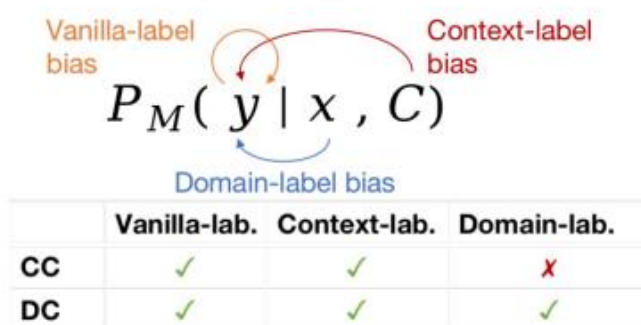
type	Precision	Recall	F1 score	Accuracy
Zeroshot	0.77	0.56	0.45	0.56
CC calibration	0.75	0.69	0.66	0.69

۳-۳- روش DC

Domain-Context Calibration

$$\bar{p}(y|C) = \frac{1}{T} \sum_{t=1}^T p(y|[Random\ i.\ d.\ text]_t, C)$$

$$\hat{y}_i = \operatorname{argmax}_{y \in \mathcal{L}} \frac{p(y|x_i, C)}{\bar{p}(y|C)}$$



Yu Fei et al., Mitigating label biases for in-context learning, ACL 2023

شکل (۳-۳) توضیح روش DC.

در این روش ابتدا ۸ کلمه رندوم از context انتخاب می‌شود و concat می‌شود و به عنوان ورودی به

مدل داده می‌شود. احتمالات خروجی دو کلاس به شرح جدول زیر است.

جدول (۳-۶) احتمالات مدل برای کلمات رندوم از context.

Positive Prob	Negative Prob
0.4234	0.5765

احتمالات فوق را به عنوان بایاس به احتمالات خروجی مدل اضافه می‌کنیم و به این صورت کالیبریشن را انجام می‌دهیم. نتایج کالیبریشن به شرح جدول زیر است.

جدول (۳-۷) مقایسه خروجی مدل برای حالت CC و zeroshot معمولی.

type	Precision	Recall	F1 score	Accuracy
Zeroshot	0.77	0.56	0.45	0.56
DC calibration	0.76	0.61	0.55	0.61

۳-۴- معیار ECE

معیار Expedted Calibration Error یک معیار برای ارزیابی کالیبریشن مدل است که طبق فرمول زیر حساب می‌شود

$$ECE = \sum \left(|Accuracy_i - Confidence_i| \times \frac{N_i}{N} \right)$$

این معیار برای دو حالت cc و dc به صورت زیر است:

type	ECE
CC	0.30
DC	0.19

نکته قابل ذکر این است چون از context کلمات را انتخاب می‌کنیم و با این احتمالات مدل را کالیبره می‌کنیم مدل بهتر می‌تواند context را تشخیص دهد.