

به نام خدا

گزارش تمرین دوم درس مدل‌های زبانی بزرگ

حمیدرضا امیرزاده

۴۰۱۲۰۶۹۹۹

## نوتبوك اول: پيش پردازش

**سوال اول:** بعد از مشاهده توكن هاى توليدشده توسط روش هاى Wordpiece و Unigram و همچنين از آنجاى كه تعداد كل توكن هاى توليدى اين دو روش تقريبا برابر است مى توان نتيجه گرفت كه توكن هاى توليدى اين دو روش تا حد زيادى مشابه هم هستند. البته روش Unigram تمايل بيشتري به جداسازى بيشتتر توكن ها دارد. به عنوان مثال توكن "بلندگو" در روش Wordpiece به همان صورت ابتدايى خود حاصل شده است اما در روش Unigram به صورت "بلند" و "گو" درآمده است.

**سوال دوم:** از آنجاى كه در اين روش سطح بايت كاراكترها درنظر گرفته شده و متن را به دنباله هاى بايت با طول متغير تجزيه مى كند، خروجى نهايى تفسيرپذيرى انساني ندارد.

**سوال سوم:** اگر سمبل Underline موجود در ابتداى توكن هاى روش Unigram را درنظر نگيريم، اشتراك خيلى زيادى بين خروجى نهايى اين دو روش وجود دارد. اما با لحاظ كردن اين سمبل، اشتراكات بيشتتر منحصر به كلماتى مى شود كه داراى پيشوند و پسوند نيستند.

## نوتبوك دوم: يادگيري در سياق

**سوال اول:** به طور كلي انتخاب معيار آريزيابي بسته به هدف آريزيابي دارد. اما از آنجايي كه در بعضي موارد پاسخ مدل درست است يا از كلمات مترادف و يا با چندين توكن اضافه تر استفاده كرده است، تكيه بر EM شايد كار درستي نباشد. به عنوان مثال داده ۱۲م داراي پاسخ اصلي Sun Life Stadium' است و مدل پاسخ Miami's Sun Life Stadium را برگردانده كه طبق معيار EM امتيازي نمي گيرد ولي جواب آن درست است.

بنابراين براي آريزيابي اين مدل روي ديتاست Squad استفاده از معيار F1 score كه حساسيت كمترى به چنين مواردى دارد و انعطاف پذير تر است پيشنهاد مي شود.

**سوال دوم:** واضح است كه پرامپت ورودى داراي اهميت زيادى در خروجى مدل است. مي توان با دادن يك پرامپت مشخص مدل را جهت دهى كرد. به عنوان مثال ذكر كنيم كه از دانش پيشين خود استفاده نكند و پاسخ هاى کوتاه بدهد. علاوه بر اين موارد، كيفيت اين پرامپت در عملکرد نهايي مدل نيز تاثير زيادى دارد.

**سوال سوم:** در حالت Answer absence از آنجايي كه هيچ راهى براي تشخيص جواب درست به كمك كانتكس هاى نامرتبط وجود ندارد، انتظار داريم كه مدل در اكثر مواقع پاسخ Not enough info را برگرداند.

در حالت Entity substitution از آنجايي كه هويت هاى موجود در كانتكس، سوال و جواب را عوض مي كنيم همچنان انتظار داريم كه مدل پاسخ مناسبى را استخراج كند. همچنين اين تغيير هويت ها باعث ايجاد متن هاى غيرواقعي مي شود كه از استفاده مدل از دانش پيشين خود جلوگیری كرده و ممكن است باعث افزايش دقت عملکرد مدل نسبت به حالت عادى شود.

در حالت Nonsense word substitution نیز به دلیل اینکه در definitions list نگاشت کلمات بی معنی به نسخه اصلی آن ها را آورده ایم، انتظار داریم مدل با استفاده از قدرت استنتاج و بازیابی خود تا حد خوبی بتواند پاسخ های درست را استخراج کند.

**سوال چهارم:** میزان دقت عملکرد مدل روی دیتاست عادی و دیتاست Answer absence می تواند نشان دهنده میزان توانایی آن در استنتاج باشد. همچنین میزان دقت عملکرد مدل روی دیتاست entity\_substitution و دیتاست Nonsense\_word\_substitution می تواند نشان دهنده میزان توانایی آن در بازیابی و همچنین استنتاج باشد.

۱- نتایج و تحلیل آن ها روی دیتاست عادی و خصمانه Answer absence:

دقت مدل روی دیتاست عادی:

EM Score=(16.93472090823084,), F1 Score=0.3671884374572664

دقت مدل روی دیتاست Answer absence (مقایسه با جواب های اصلی دیتاست):

EM Score=(0.9460737937559129,), F1 Score=0.01965991172715234

همانطور که مشاهده می شود، افت عملکرد شدیدی در این حالت وجود دارد که علت آن در بخش قبل توضیح داده شد.

دقت مدل روی دیتاست Answer absence (مقایسه با Not enough info):

EM Score=80.79470198675497

از آنجایی که در اینجا دقیقاً میخواستیم خروجی مدل همان Not enough info باشد، تنها از معیار EM استفاده کردیم. مشاهده می شود که در این حالت مدل عملکرد خوبی در اظهار نادانی خود داشته است.

در بحث تحلیل پاسخ های مدل می توان به مثال ۲۰ام دیتاست اشاره کرد. جایی که پاسخ درست برابر the Panthers بوده است ولی مدل پاسخ the Dallas Cowboys را برگردانده. در اینجا باید مدل به دلیل نداشتن کانتکس لازم باید جواب Not enough info را برمی گرداند ولی با استفاده از اطلاعات موجود در وزن های خود پاسخ غلط گفته شده را توهم زده است.

اما به طور کلی از آنجایی دقت EM مدل با مقایسه خروجی آن با Not enough info برابر ۸۰٪ است، می توان اینگونه نتیجه گیری کرد مدل اتکای نسبتا خوبی به کانتکس ورودی خود دارد و همانطور که در پرامپت گفته شده، اگر جواب سوالی را نمی داند، پاسخ نمی دهد. در سایر موارد مشخصا مدل توهم زده و پاسخ های غلط مربوط به کانتکس نامرتبط ورودی خود برمی گرداند. با کمک بررسی های دستی می توان مشاهده کرد که در حالت هایی که مدل پاسخ می دهد، بایاس به این سمت دارد که اسم خاص برگرداند.

## ۲- نتایج و تحلیل آن ها روی دیتاست عادی و خصمانه Entity Substitution:

دقت مدل روی دیتاست عادی:

EM Score=(16.93472090823084,), F1 Score=0.3671884374572664

دقت مدل روی دیتاست Entity Substitution:

EM Score=(21.665089877010406,), F1 Score=0.3788809504407445

مشاهده می شود در این حالت افت عملکرد اتفاق نیفتاده بلکه بهبود نیز داشته است. دلیل آن را می توان غیرواقعی شدن سناریوها و عدم توانایی مدل در ساختن پاسخ های مبتنی بر دانش پیشین خود دانست.

در بحث تحلیل پاسخ های مدل می توان به یک مثال جالب اشاره کرد. مثال ۷۸ام دیتاست پاسخ درست Anderson بوده که مدل پاسخ نادرست Newton را برگردانده، اما بعد از عوض کردن هویت های نامدار پاسخ درست به Christiano ronaldo تغییر کرده و مدل نیز پاسخ درست

ronaldo را برگردانده و به نوعی اشتباه خود در ابتدا را تکرار نکرده است. دلیل وجود چنین مواردی همانطور که اشاره شد، عدم توانایی مدل در توهّم زدن پاسخ در سناریوهای غیرواقعی است.

### ۳- نتایج و تحلیل آن ها روی دیتاست عادی و خصمانه Nonsense Word Substitution:

دقت مدل روی دیتاست عادی:

EM Score=(16.93472090823084,,), F1 Score=0.3671884374572664

دقت مدل روی دیتاست Nonsense Word Substitution:

EM Score=(3.5004730368968775,,), F1 Score=0.20720127477626685

لازم به ذکر است که در اینجا در پرامپت ورودی به مدل ذکر شده است که به عنوان پاسخ سعی کند به کمک definitions list کلمه اصلی را به جای کلمه بی معنی برگرداند. مشاهده می شود که از نظر معیار EM نسبت به حالت عادی عملکرد به طور قابل ملاحظه ای افت کرده است اما همچنان خیلی بهتر از حالت Answer absence است و این نشان دهنده وجود قابلیت بازیابی و اسنتاج در مدل است که باعث می شود بتواند تا حدی عمیات استخراج جواب و نگاشت آن به کلمه اصلی خود را انجام دهد.

**سوال پنجم:** در حالت های Answer absence و Entity Substitution همانگونه که در ابتدا توضیح داده شد و نتایج مشاهده شد، عملکرد مدل مطابق انتظارات قبلی بود. اما در مورد Nonsense Word Substitution با توجه به تفهیم مدل به اینکه باید کلمه اصلی را برگرداند اما همچنان دقت نهایی آن با دقت حالت عادی تفاوت زیادی داشت و انتظار این بود که عملکرد مناسبتری داشته باشد. این نشان دهنده قدرت بازیابی و اسنتاج نسبتا پایین مدل Llama2- 7B می باشد.

**سوال ششم:** با استفاده از طراحی نمونه‌های خصمانه هدفمند می‌توان توانایی مدل در موارد خاص مثل استنتاج و یا بازیابی را هدف قرار داد. به این صورت که اگر مدل بتواند حتی با حضور عوامل گمراه کننده موجود در نمونه های خصمانه به پاسخ درست برسد، نشان دهنده مقاومت مدل است. به طور خاص با قرار دادن یک کانتکس نامرتبط توانایی مدل در استنتاج اینکه پاسخ در کانتکس وجود ندارد را می‌توان ارزیابی کرد. و یا با تعویض هویت های نامدار و ساخت کلمات مصنوعی می‌توان توانایی مدل در بازیابی اطلاعات را ارزیابی کرد. به طور کلی می‌توان گفت که این ارزیابی‌ها به درک عمیق‌تر محدودیت‌ها و توانایی های خاص مدل کمک می‌کند.

## نوتبوك سوم: تنظيم سازي

### • Zero-shot / Few-shot setting

در جدول زير نتايج روش هاي بدون مثال و چندمثال آورده شده است.

Approach	Pos0	Pos1	Pos2	Neg0	Neg1	Neg2	F1-score (0)	F1-score (1)	F1-score (Acc)
Zero-shot with 'positive' and 'negative' labels	-	-	-	-	-	-	0.21	0.69	0.56
Zero-shot with '0' and '1' labels	-	-	-	-	-	-	0.14	0.67	0.52
Few-shot	✓	-	-	✓	-	-	0.64	0.78	0.73
Few-shot	✓	-	-	-	✓	-	0.57	0.76	0.69
Few-shot	✓	-	-	-	-	✓	0.76	0.81	0.79
Few-shot	-	✓	-	✓	-	-	0.52	0.74	0.67
Few-shot	-	✓	-	-	✓	-	0.54	0.75	0.67
Few-shot	-	✓	-	-	-	✓	0.71	0.79	0.75
Few-shot	-	-	✓	✓	-	-	0.58	0.76	0.69
Few-shot	-	-	✓	-	✓	-	0.40	0.72	0.62
Few-shot	-	-	✓	-	-	✓	0.57	0.75	0.69
Few-shot	1	2	-	3	-	-	0.80	0.76	0.78
Few-shot	1	3	-	2	-	-	0.72	0.79	0.76
Few-shot	2	1	-	3	-	-	0.81	<b>0.83</b>	<b>0.82</b>
Few-shot	3	1	-	2	-	-	<b>0.83</b>	0.81	<b>0.82</b>
Few-shot	2	3	-	1	-	-	0.79	0.70	0.75
Few-shot	3	2	-	1	-	-	0.73	0.43	0.63



باتوجه به نتایج آورده شده همانطور که انتظار می‌رفت دقت روش چندمثال بهتر از روش بدون مثال می‌باشد، همچنین عوض کردن لیبل خروجی در حالت بدون مثال نیز می‌تواند بر عملکرد مدل تاثیرگذار باشد.

در حالت چندمثال مشاهده می‌شود بسته به اینکه کدام یک از مثال‌ها و با چه ترتیبی در پرامپت ورودی به مدل داده شوند، روی دقت مدل تاثیر می‌گذارد و بنابراین خروجی مدل به این عوامل بستگی زیادی دارد. توجه شود که در حالت چندمثال آخر از آنجایی که یک مثال اضافه تر نسبت به حالت چندمثال قبلی داده می‌شود، عموماً نتیجه بهتری حاصل می‌شود. علاوه بر این معمولاً وقتی مثال آخر دارای لیبل ۰ باشد، دقت مدل در پیشبینی داده‌های ۰ نیز بهتر می‌شود و برعکس.

## • Calibration setting

در جدول زیر نتایج روش های بدون مثال کالیبره شده با استفاده از روش های CC و DC آورده شده است.

با مقایسه نتایج این روش های تنظیم شده نسبت به روش غیر تنظیم شده در حالت بدون مثال می‌توان اثر مثبت تنظیم سازی را مشاهده کرد. همچنین روش DC نسبت به روش CC عملکرد بهتری دارد زیرا زمانی که از توکن های مطابق با سیاق تسک برای تنظیم کردن مدل استفاده کنیم، به تخمین بهتری برای احتمالات کالیبره سازی نسبت به حالتی که فقط از توکن از N/A برای این کار استفاده کنیم می‌رسیم.

Approach	neg_prob_calibration	pos_prob_calibration	F1-score (0)	F1-score (1)	F1-score (Acc)	ECE
<b>CC calibrated Zero-shot</b>	0.865	0.135	0.75	0.53	0.67	0.0502
<b>DC calibrated Zero-shot</b>	0.805	0.195	<b>0.76</b>	<b>0.70</b>	<b>0.74</b>	0.0883

