

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی اصفهان

دانشکده مهندسی برق و کامپیوتر

گزارش پروژه‌ی کارشناسی

عنوان پروژه:

سیستم پرسش و پاسخ با استفاده از یادگیری ماشین

استاد راهنما:

دکتر مهران صفایانی

توسط:

مرضیه نوری

بهمن ۱۴۰۰

فهرست مطالب

چکیده	۱
فصل اول: مقدمه	۲
فصل دوم: معرفی داده‌های پروژه	۴
۱-۲ معرفی مجموعه داده‌های مرتبط با پرسش و پاسخ	۴
۱-۱-۲ SQuAD	۴
۲-۱-۲ PersianQA	۷
۲-۲ ملاک‌های ارزیابی	۹
فصل سوم: روش‌های متداول پردازش زبان طبیعی	۱۱
۱-۳ شبکه‌های عصبی بازگشتی	۱۱
۲-۳ Long Short Term Memory	۱۳
فصل چهارم: روش‌های مبتنی بر مبدل	۱۶
۱-۴ مبدل	۱۶
۱-۱-۴ معماری مبدل	۱۷
۲-۱-۴ رمزگذاری موقعیتی	۲۵
۲-۴ مدل BERT	۲۶
۱-۲-۴ معماری BERT	۲۶
۲-۲-۴ نمایش ورودی در مدل BERT	۲۷
۳-۲-۴ رویکردهای استفاده از مدل BERT	۲۸
۴-۲-۴ استفاده از مدل BERT با رویکرد تنظیم دقیق	۲۸
۳-۴ مدل ParsBERT	۳۱

۳۱	۴-۳-۱ جمع‌آوری داده
۳۱	۴-۳-۲ پیش‌پردازش داده
۳۲	۴-۳-۳ تقسیم اسناد به جملات صحیح فارسی
۳۴	فصل پنجم: حل مسئله
۳۴	۵-۱ تنظیم دقیق ParsBERT برای کار پرسش و پاسخ
۳۶	۵-۲ پیاده‌سازی
۳۶	۵-۲-۱ آماده‌سازی داده‌ها
۳۶	۵-۲-۲ تنظیم دقیق مدل
۳۷	۵-۲-۳ ارزیابی مدل
۳۸	۵-۳ نتایج ارزیابی
۴۱	نتیجه‌گیری
۴۲	مراجع

چکیده

امروزه کم تر پیش می آید در روز حداقل یک بار برای یافتن پاسخ سوالی که داریم، به موتورهای جستجویی نظیر گوگل متوسل نشویم. آمار نشان می دهد روزانه ۵/۶ میلیارد جستجو روی گوگل انجام می شود و این تنها مربوط به یکی از چندین موتور جستجوی مطرح است. با افزایش سرعت تولید داده ها، از جمله داده های متنی، نیاز به سیستم هایی که علاوه بر استخراج اسناد مرتبط با سوال ما، پاسخ مستقیم سوالمان را نیز به ما برگردانند، بیشتر شده است. در این پروژه، برای آنکه قدمی در راستای حل این چالش برداریم، مدلی طراحی کرده ایم که مبتنی بر مبدل ParsBERT^۱ بوده و روی مجموعه داده ی PersianQA^۲ که یک مجموعه ی پرسش و پاسخ فارسی است، تنظیم دقیق^۴ شده است. در این پروژه نهایتاً با تنظیم ابرپارامترها و آزمایش یازده حالت، به بالاترین امتیاز F1^۵ ۵۹,۱۲٪ و امتیاز تطابق دقیق^۶ (EM) ۴۲,۷۹٪ رسیده ایم.

^۱ <https://blog.hubspot.com/marketing/google-search-statistics>

^۲ transformer

^۳ Persian Question Answering

^۴ fine-tune

^۵ F1 score

^۶ exact match (EM)

فصل اول

مقدمه

امروزه با توجه به رشد روزافزون روند تولید داده‌ها، اهمیت آنالیز و بهره‌جستن از آن‌ها به کمک یادگیری ماشین، بر کسی پوشیده نیست؛ یکی از زیرشاخه‌های مهم این حوزه، پردازش زبان‌های طبیعی است که خود شاخه‌ای وسیع به شمار می‌آید. هدف در پردازش زبان‌های طبیعی آنست که سیستم‌های کامپیوتری بتوانند زبان‌های طبیعی مانند انگلیسی و فارسی را درک کنند و در اموری مانند طراحی سیستم‌های پرسش و پاسخ^۲، ترجمه‌ی ماشینی^۳، تحلیل احساس^۴، تشخیص گفتار^۵، خلاصه‌سازی اسناد^۶ و ... به کمک ما بیایند.

^۱ Natural Language Processing (NLP)

^۲ Question Answering Systems

^۳ Machine Translation

^۴ Sentiment Analysis

^۵ Speech Recognition

^۶ Document Summarization

همانطور که اشاره شد، یکی از کاربردهای پردازش زبان‌های طبیعی، بحث طراحی سیستم‌های پرسش و پاسخ است. در دنیایی که روزانه توده‌ی عظیمی از اطلاعات تولید می‌شود، به نظر می‌رسد آنچه که در حالت عادی، موتورهای جستجو در دسترس ما قرار می‌دهند، یعنی اسنادی مرتبط با سوالات ما، به اندازه‌ی گذشته کارآمد نیست؛ چرا که هدف اغلب افراد از جستجو بر بستر وب، یافتن پاسخی مشخص برای یک پرسش است و ترجیح بر آنست که مستقیماً بتوانند پاسخ سوال خود را دریافت کنند؛ تا آنکه سندی مرتبط با موضوع مورد پرسش خود ببینند که احتمالاً پاسخ سوالشان را در خود دارد. خصوصاً با گسترش استفاده از تلفن همراه و جستجو توسط آن، افراد عموماً تمایلی به خواندن متن‌های بلندبالا که فقط بخشی از آن برایشان مفید خواهد بود، ندارند.

بر این اساس در این پروژه، می‌خواهیم به کمک یادگیری ماشین، سیستم پرسش و پاسخی به زبان فارسی طراحی کنیم که با داشتن یک متن و یک پرسش، بتواند پاسخ را از متن استخراج کند. این سیستم‌ها عموماً با عنوان سیستم‌های درک مطلب^۱ شناخته می‌شوند. سیستم طراحی شده، مبتنی بر مبدل^۲ بوده و با استفاده از مدل ParsBERT، و به کمک مجموعه داده‌ی PersianQA برای این کار تنظیم دقیق شده است.

^۱ reading comprehension systems

^۲ transformer

^۳ fine-tune

فصل دوم

معرفی داده‌های پروژه

در این فصل به معرفی مجموعه داده‌های مرتبط با پرسش و پاسخ خواهیم پرداخت و سپس ملاک‌های ارزیابی سیستم‌هایی که با این مجموعه داده‌ها سازگار می‌شوند را معرفی خواهیم کرد.

۱-۲ معرفی مجموعه داده‌های مرتبط با پرسش و پاسخ

در این پروژه از مجموعه داده‌ی PersianQA^۱ استفاده کرده‌ایم. این مجموعه داده، از مجموعه‌ی SQuAD^۲ الهام گرفته است. در این بخش به معرفی این دو مجموعه می‌پردازیم.

SQuAD ۱-۱-۲

SQuAD یک مجموعه داده‌ی درک مطلب به زبان انگلیسی است. پاراگراف‌های این مجموعه داده، برگرفته از مقالات

^۱ Persian Question Answering

^۲ Stanford Question Answering Dataset

ویکی‌پدیا بوده و پرسش‌ها و پاسخ‌های آن، جمع‌سپاری^۱ شده‌اند. این مجموعه در دو نسخه ارائه شده است:

- SQuAD 1.1 شامل بیش از صد هزار پرسش است که از بیش از پانصد مقاله‌ی ویکی‌پدیا استخراج شده‌اند. تمام این سوالات را می‌توان با توجه به متن پاسخ داد. پاسخ به هر سوال، گستره‌ای از متن را شامل می‌شود. [۱۳]
- SQuAD 2.0 علاوه بر صد هزار پرسش نسخه‌ی پیشین، شامل بالغ بر پنجاه هزار سوال بدون جواب نیز هست. در واقع پاسخ به این سوالات، با استفاده از متن، ممکن نیست. از این سوالات، با عنوان سوالات بدون جواب یاد می‌کنیم. [۱۲]

نکته‌ی مهم در رابطه با این مجموعه داده آنست که اگر پرسشی قابل پاسخ‌دهی باشد، پاسخ آن، قطعا گستره‌ای از متن است. این بدان معناست که سیستم‌هایی که از این مجموعه داده استفاده می‌کنند، متن پاسخ را تولید نمی‌کنند؛ بلکه آن را از مقاله‌ای که در اختیار دارند استخراج می‌سازند؛ یعنی بازه‌ای به هم پیوسته از پاراگراف را به عنوان پاسخ، انتخاب می‌کنند. این کار شبیه زمانبست که بخشی از یک متن را برجسته^۳ می‌کنیم.

هر نمونه از SQuAD از سه بخش اصلی پرسش، متن و پاسخ تشکیل شده است. در زیر، نمونه‌ای از یک پرسش با پاسخ از این مجموعه آورده شده است:

answer	context	question	title
{'answer_start': [727], 'text': ['animal electricity']}	William Champion's brother, John, patented a process in 1758 for calcining zinc sulfide into an oxide usable in the retort process. Prior to this, only calamine could be used to produce zinc. In 1798, Johann Christian Ruberg improved on the smelting process by building the first horizontal retort smelter. Jean-Jacques Daniel Dony built a different kind of horizontal zinc smelter in Belgium, which processed even more zinc. Italian doctor Luigi Galvani discovered in 1780 that connecting the spinal cord of a freshly dissected frog to an iron rail attached by a brass hook caused the frog's leg to twitch. He incorrectly thought he had discovered an ability of nerves and muscles to create electricity and called the effect "animal electricity". The galvanic cell and the process	What did Galvani name the effect he created of causing the frogs legs to twitch?	Zinc

¹ crowdsourced

² span

³ highlight

	of galvanization were both named for Luigi Galvani and these discoveries paved the way for electrical batteries, galvanization and cathodic protection.		
--	---	--	--

همانطور که در نمونه‌ی بالا می‌بینیم، هر پرسش قابل پاسخی، علاوه بر متن جواب، شاخص^۱ شروع جواب از متن را نیز در بردارد که برای مثال آورده شده، این عدد برابر با ۷۲۷ است. این شاخص نشان می‌دهد پاسخ پرسش ما از چندمین حرف مقاله شروع می‌شود. واضح است که شاخص پایان جواب از مجموع شاخص شروع جواب و طول جواب به دست می‌آید.

نمونه‌ای از یک پرسش بدون جواب از این مجموعه را در زیر می‌بینیم:

anwer	context	question	title
<p>text: [], '}</p> <p>{[]:'answer_start</p>	<p>In addition to Old Persian and Avestan, which are the only directly attested Old Iranian languages, all Middle Iranian languages must have had a predecessor "Old Iranian" form of that language, and thus can all be said to have had an (at least hypothetical) "Old" form. Such hypothetical Old Iranian languages include Carduchi (the hypothetical predecessor to Kurdish) and Old Parthian. Additionally, the existence of unattested languages can sometimes be inferred from the impact they had on neighbouring languages. Such transfer is known to have occurred for Old Persian, which has (what is called) a "Median" substrate in some of its vocabulary. Also, foreign references to languages can also provide a hint to the existence of otherwise unattested languages, for example through toponyms/ethnonyms or in the recording of vocabulary, as Herodotus did for what he called</p> <p>"Scythian</p>	<p>What cannot be inferred from the impact on a neighboring ?language</p>	<p>Iranian_languages</p>

¹ index

همانطور که مشاهده می‌شود، برای پاسخ به این سوال که «چه چیزی را نمی‌توان از تأثیر بر زبان همسایه استنباط کرد؟» در متن داده‌شده، منبع کافی وجود ندارد و این سوال بدون جواب طبقه‌بندی می‌شود.

SQuAD علاوه بر مجموعه‌ی آموزش، شامل دو مجموعه‌ی دیگر مربوط به اعتبارسنجی^۱ و آزمایش^۲ نیز هست. مجموعه‌ی آموزش و اعتبارسنجی به شکل عمومی منتشر شده‌اند اما مجموعه‌ی آزمایش به صورت عمومی در دسترس نیست. در مجموعه‌ی آزمایش و اعتبارسنجی، هر پرسش، دارای سه پاسخ است که لزوماً عین یکدیگر نیستند اما هر سه پاسخ درستی ارزیابی می‌شوند. دلیل این امر آنست که هر یک از این پرسش‌ها توسط شخص متفاوتی پاسخ داده شده‌اند. این امر با هدف مقاوم‌کردن ارزیابی صورت گرفته است.

PersianQA ۲-۱-۲

PersianQA یک مجموعه‌داده‌ی درک مطلب به زبان فارسی است. هر نمونه از این مجموعه‌داده نیز مشابه SQuAD از سه جزء اساسی تشکیل شده است که شامل پرسش، متن و پاسخ می‌باشد. متن‌های این مجموعه‌داده، از روی مقالات ویکی‌پدیای فارسی و از دسته‌بندی‌های موضوعی مختلفی همچون تاریخی، جغرافیایی، مذهبی، علمی و ... استخراج شده‌اند. پرسش و پاسخ‌های مربوط به این مقالات، توسط افراد مختلفی که فارسی، زبان مادری آن‌هاست، صورت گرفته است. لازم به ذکر است بعضی از پرسش‌ها به زبان محاوره‌ای آورده شده‌اند. [۲]

این مجموعه‌داده شامل بیش از ۹۰۰۰ پرسش است که مانند نسخه‌ی دوم SQuAD، دارای سوالات بدون جواب نیز می‌باشد. علاوه بر این ۹۰۰۰ پرسش، حدود ۹۰۰ پرسش دیگر نیز به عنوان دادگان آزمایش وجود دارد. به طور متوسط، هر پاراگراف، هفت پرسش قابل پاسخ و سه پرسش غیرقابل پاسخ دارد. هر یک از پرسش‌های موجود در نمونه‌های آزمایش، دارای دو پاسخ هستند که همچون SQuAD لزوماً برابر یکدیگر نیستند.

در زیر نمونه‌ای از یک پرسش دارای جواب از این مجموعه آورده شده است:

پاسخ	پرسش	متن	موضوع
{'text': ['صابون، به تعریف علم شیمی، نمک	صابون چییه؟	نخستین اشاره به ساختن یک ماده شبیه صابون مربوط به حدود ۲۸۰۰ سال پیش از میلاد در بابل باستان است. برخی منابع مصری مربوط به حدود ۱۵۰۰ پیش از میلاد نیز به ساختن چنین ماده‌ای اشاره دارند. در	صابون

^۱ validation

^۲ test

^۳ robust

حدود ۶۰۰ سال پیش از میلاد، ملوانان فینیقی‌ها، فن صابون‌سازی (یا صابون پزی) را به سواحل مدیترانه بردند. در قرن اول میلادی، بهترین صابون از چربی بز و خاکسترهای به‌دست‌آمده از سوزاندن چوب «درخت آتش» به‌دست می‌آمد. تا پایان قرن هجدهم، صابون را از چربی حیوانی و خاکستر چوب تهیه می‌کردند. در همان هنگام، معلوم شد که می‌توان به‌جای خاکستر چوب، از سود سوزآور، که قلیای حاصل از نمک معمولی است، استفاده کرد. صابون، به تعریف علم شیمی، نمک یک اسید چرب است. صابون عمدتاً برای شست‌وشو، حمام کردن و پاکیزگی استفاده می‌شود، ولی در رسیدن پارچه هم از صابون استفاده می‌گردد و جزء مهمی از روان‌سازها است. صابون توالست سخت با بوی دلپذیر در خاورمیانه در دوره طلایی اسلامی تولید شد، زمانی که صابون سازی به یک صنعت تبدیل شد. دستور العمل‌های تهیه صابون توسط محمد بن زکریا رازی شرح داده شده‌است، که همچنین دستور تهیه گلیسرین از روغن زیتون را داده است.	یک اسید چرب است. , 'answer_start': [585]]
---	--

مشابه آنچه در SQuAD دیدیم، هر پرسش دارای پاسخ، علاوه بر متن جواب، شاخص شروع جواب از متن را نیز در بردارد که برای مثال آورده شده این عدد برابر با ۵۸۵ است.

نمونه‌ای از یک پرسش بدون جواب از این مجموعه در زیر آورده شده است:

پاسخ	پرسش	متن	موضوع
{'text': [], 'answer_start': []}	اوج گیری تصوف چه اثری روی ابوبکر بن سعد داشت؟	ابومحمد مُشرف‌الدین مُصلِح بن عبدالله بن مشرّف (۶۰۶ – ۶۹۰ هجری قمری) متخلص به سعدی، شاعر و نویسنده پارسی‌گوی ایرانی است. اهل ادب به او لقب «استادِ سخن»، «پادشاهِ سخن»، «شیخِ اجل» و حتی به‌طور مطلق، «استاد» داده‌اند. او در نظامیه بغداد، که مهم‌ترین مرکز علم و دانش جهان اسلام در آن زمان به حساب می‌آمد، تحصیل و پس از آن به‌عنوان خطیب به مناطق مختلفی از جمله شام و حجاز سفر کرد. سعدی سپس به زادگاه خود شیراز، برگشت و تا پایان عمر در آن‌جا اقامت گزید. آرامگاه سعدی در شیراز واقع شده‌است که به سعدیه معروف است. بیشتر عمر او مصادف با حکومت اتابکان فارس در شیراز	سعدی

	و هم‌زمان با حمله مغول به ایران و سقوط بسیاری از حکومت‌های وقت نظیر خوارزمشاهیان و عباسیان بود. البته سرزمین فارس، به واسطه تدابیر ابوبکر بن سعد، ششمین و معروف‌ترین اتابکان سلجوقی شیراز، از حمله مغول در امان ماند. همچنین قرن ششم و هفتم هجری مصادف با اوج‌گیری تصوف در ایران بود و تأثیر این جریان فکری و فرهنگی در آثار سعدی قابل ملاحظه است.	
--	---	--

۲-۲ ملاک‌های ارزیابی

عملکرد سیستمی که با استفاده از PersianQA خواهیم ساخت، با دو ملاک امتیاز تطابق دقیق^۱ و امتیاز F1^۲ ارزیابی می‌شود.

- تطابق دقیق یک معیار باینری (درست/نادرست) است که نشان می‌دهد آیا خروجی سیستم دقیقاً با پاسخ درست مطابقت دارد یا خیر. به عنوان مثال، اگر سیستم ما، «انیشتین» را به عنوان پاسخ یک سوال برگرداند، در حالی که پاسخ صحیح «آلبرت انیشتین» باشد، در آن صورت سیستم نمره‌ی تطابق دقیق صفر برای آن مثال، دریافت خواهد کرد. تطابق دقیق یک معیار سخت‌گیرانه است.

- F1 یک ملاک کمتر سختگیرانه است. برای این معیار، دو پارامتر مطرح است: دقت^۳ و یادآوری^۴. در مثال «انیشتین»، سیستم دارای دقت (precision) ۱۰۰٪ است؛ چرا که پاسخ آن یعنی انیشتین، زیرمجموعه‌ای از پاسخ درست یعنی آلبرت انیشتین می‌باشد. پارامتر یادآوری (recall) برای این مثال برابر با ۵۰٪ است؛ چرا که فقط شامل یکی از دو کلمه‌ایست که در پاسخ اصلی ذکر شده است. امتیاز F1 از طریق رابطه‌ی (۱-۲) به دست می‌آید:

$$F1 = 2 \times \frac{precision \times recall}{precision + recall} \quad (1-2)$$

بنابراین در مثال بالا، این عدد برابر خواهد بود با:

$$F1 = 2 \times \frac{100 \times 50}{100 + 50} = 66.67\%$$

^۱ Exact Match (EM)

^۲ F1 score

^۳ precision

^۴ recall

هنگام ارزیابی روی مجموعه‌های اعتبارسنجی یا آزمایش، حداکثر امتیازات $F1$ و تطابق دقیق را در دو پاسخ ارائه‌شده توسط انسان برای آن سؤال در نظر می‌گیریم؛ به عنوان مثال، اگر سوالی که پاسخ آن آلبرت انیشتین بوده در مجموعه‌ی اعتبارسنجی حضور داشته باشد و پاسخ اول «آلبرت انیشتین» و پاسخ دوم «انیشتین» باشد، سیستم برای این مثال، تطابق دقیق ۱۰۰ درصد و امتیاز $F1$ ۱۰۰ درصد را دریافت می‌کند. در نهایت، نمرات تطابق دقیق و $F1$ در کل مجموعه داده‌ی ارزیابی برای به دست آوردن نمرات گزارش شده‌ی نهایی به طور میانگین محاسبه می‌شوند.

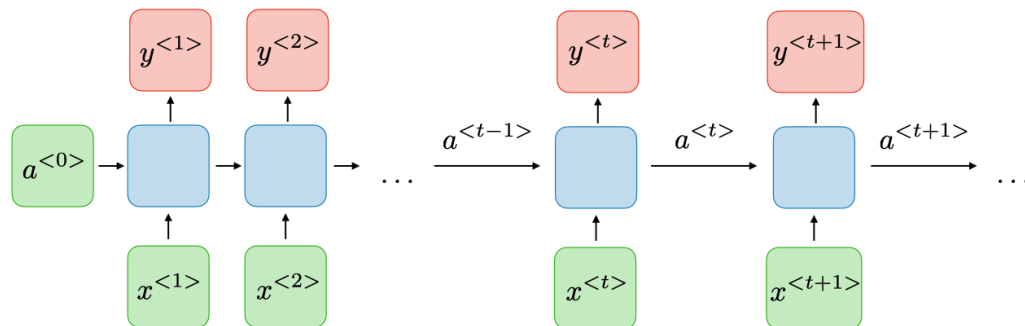
فصل سوم

روش‌های متداول پردازش زبان طبیعی

در این فصل، با روش‌های مرسوم که در پردازش زبان طبیعی به کار می‌روند، آشنا می‌شویم. نظر به اینکه پروژه‌ی پیاده‌سازی‌شده در این پروژه، مبتنی بر مبدل‌هاست، از ذکر جزئیات بیش از حد در این فصل پرهیز کرده‌ایم.

۱-۳ شبکه‌های عصبی بازگشتی

شبکه‌های عصبی بازگشتی که با نام RNN نیز شناخته می‌شوند، دسته‌ای از شبکه‌های عصبی هستند که در هر مرحله، علاوه بر ورودی جدید، از خروجی مرحله‌ی قبل نیز تاثیر می‌گیرند. در واقع خروجی تولید شده در هر گام از زمان، تنها متکی به ورودی جدید نیست و از خروجی گام‌های پیشین نیز متاثر است. شکل ۱-۳ معماری یک مدل ساده از شبکه‌ی عصبی بازگشتی را نشان می‌دهد.



شکل (۱-۳) معماری شبکه‌های عصبی بازگشتی [۱]

همانطور که می‌بینیم، در گام t ام، ورودی جدید $x^{(t)}$ به همراه $a^{(t-1)}$ که خروجی گام قبل است، وارد شبکه می‌شود و خروجی $y^{(t)}$ را تولید می‌کند.

ذات ترتیبی این شبکه‌ها، آن‌ها را به شبکه‌های مناسبی برای داده‌های ترتیبی، مانند متن و گفتار تبدیل کرده است.

در شبکه‌ای آورده شده، به ازای هر ورودی، یک خروجی دریافت می‌کنیم. به این معماری، شبکه‌ای بازگشتی چند به چند^۱ با تعداد ورودی و خروجی برابر گفته می‌شود. شبکه‌ای چند به چند یعنی ما از چند ورودی، به چند خروجی می‌رسیم. از این معماری می‌توان برای کارهایی هم‌چون تشخیص موجودیت‌های نام‌دار^۲ استفاده کرد؛ اما باید گفت این مدل، تنها ساختاری نیست که می‌توان شبکه‌های عصبی بازگشتی را به کار برد. به فراخور هدفی که دنبال می‌کنیم می‌توانیم از دیگر معماری‌های بازگشتی نیز استفاده کنیم. از این معماری‌ها می‌توان به معماری یک به چند^۳ (مورد استفاده در تولید موسیقی^۴)، معماری چند به یک^۵ (مورد استفاده در طبقه‌بندی احساس^۶) و معماری چند به چند با تعداد ورودی و خروجی نابرابر (مورد استفاده در ترجمه ماشینی) اشاره کرد.

از مزایای این شبکه‌ها می‌توان به موارد زیر اشاره کرد:

- برخلاف شبکه‌های عصبی استاندارد، در شبکه‌های عصبی بازگشتی، پردازش ورودی با هر طولی ممکن است و نیاز نیست ورودی حتما دارای طول مشخصی باشد.

^۱ many to many

^۲ Named Entity Recognition (NER)

^۳ one to many

^۴ music generation

^۵ many to one

^۶ sentiment classification

- با توجه به تاثیری که این شبکه از گذشته می گیرد، اطلاعات در طول زمان منتقل می شوند.
- اندازه‌ی مدل با افزایش اندازه‌ی ورودی تغییر نمی کند و فقط به تعداد گام‌های زمانی بیشتری برای پردازش نیاز است.

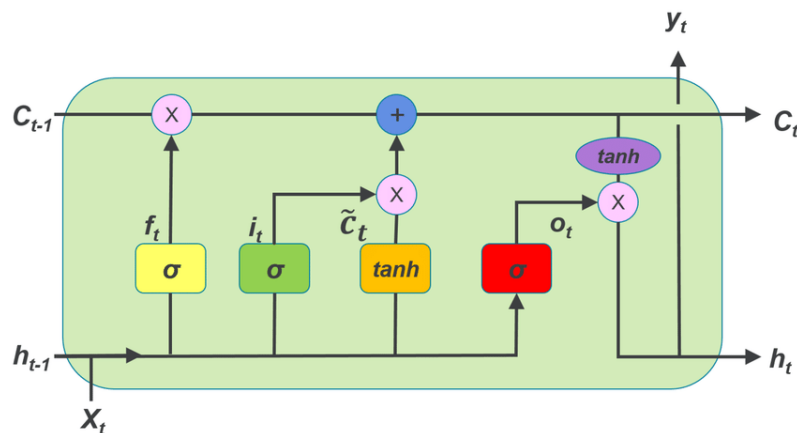
در کنار مزایایی که برای شبکه‌های عصبی بازگشتی، برشمردیم، تعدادی ایراد جدی نیز متوجه آن‌ها است:

- به این دلیل که پردازش ورودی باید به صورت متوالی انجام گیرد، این شبکه‌ها به کندی آموزش می بینند و پردازش موازی در آن‌ها ممکن نیست.
- این نوع شبکه‌ها با اینکه از گذشته تاثیر می گیرند، اما در طول زمان، وابستگی‌های بلندمدت را فراموش می کنند.
- این معماری، با مشکل محو و گرادیان روبه‌روست.

برای رفع دو مشکل محو گرادیان و فراموشی وابستگی‌ها، نوعی از شبکه‌های عصبی بازگشتی با نام Long Short Term Memory که به اختصار از آن‌ها با عنوان LSTM یاد می شود، معرفی شد که در بخش بعد به آن می پردازیم.

۲-۳ Long Short Term Memory

شبکه‌ی LSTM با هدف حل مشکل محو گرادیان و هم چنین حفظ طولانی تر وابستگی‌ها نسبت به RNN ساده، ایجاد شد. [۵] در این مدل، سه دروازه و یک سلول حافظه (c) معرفی شدند. این دروازه‌ها عبارتند از دروازه‌ی ورودی^۱ (i)، دروازه‌ی خروجی^۲ (o) و دروازه‌ی فراموشی^۳ (f). شکل ۲-۳ معماری داخلی این شبکه را به صورت شماتیک نشان می دهد.



شکل (۲-۳) ساختار یک بلوک LSTM [۸]

¹ input gate

² output gate

³ forget gate

دروازه‌ی i مشخص می‌کند چه مقدار از اطلاعات فعلی (گام t) یعنی \tilde{C}_t در حافظه ذخیره شده و وارد مرحله‌ی بعد شود. برای به دست آوردن \tilde{C}_t باید خروجی گام قبلی یعنی h_{t-1} و ورودی جدید یعنی x_t در کنار هم قرار بگیرند. سپس در ماتریس وزن W_C که در طول فرآیند آموزش مقادیر آن به دست آمده است، ضرب شود. بعد از آن با بایاس^۱ b_C که آن هم در طول یادگیری، به دست آمده است جمع گردد و نهایتاً از یک تابع فعالسازی \tanh عبور کند. برای محاسبه‌ی دروازه‌ی i در گام فعلی، به شکل مشابه، باید خروجی گام قبلی یعنی h_{t-1} و ورودی جدید یعنی x_t در کنار هم قرار گرفته و در ماتریس وزن W_i که همراه با مدل آموزش دیده است، ضرب شود. بعد با بایاس b_i که آن هم در طول یادگیری، به دست آمده است جمع شده و در آخر از یک تابع فعالسازی sigmoid عبور کند. صورت ریاضی این دو عملیات به ترتیب در فرمول‌های (۱-۳) و (۲-۳) آمده است.

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i) \quad (1-3)$$

$$\tilde{C}_t = \tanh(W_C[h_{t-1}, x_t] + b_C) \quad (2-3)$$

دروازه‌ی f تعیین می‌کند چه میزان از اطلاعات گذشته از سلول حافظه پاک شود. برای به دست آوردن خروجی این دروازه در گام t ، باید خروجی گام قبلی یعنی h_{t-1} و ورودی جدید یعنی x_t الحاق شوند و پس از آن در ماتریس وزن W_f که همراه با مدل آموزش دیده است، ضرب شود. بعد با بایاس b_f که آن هم در طول یادگیری، به دست آمده جمع گردد و از یک تابع فعالسازی \tanh عبور کند. صورت ریاضی این عملیات در فرمول (۳-۳) آمده است.

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f) \quad (3-3)$$

نهایتاً C_t با رابطه‌ی (۴-۳) محاسبه می‌شود. دقت داریم که در این فرمول علامت $*$ به معنای ضرب عنصر به عنصر است.

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \quad (4-3)$$

دروازه‌ی o مشخص می‌کند چه اطلاعاتی از حافظه خوانده شده و به خروجی انتقال یابد. برای محاسبه‌ی این دروازه، عملیاتی مشابه دو دروازه‌ی قبلی صورت می‌گیرد. یعنی باید خروجی گام قبلی یعنی h_{t-1} و ورودی جدید یعنی x_t در کنار هم قرار بگیرند. سپس در ماتریس وزن W_o که آموزش دیده است، ضرب شود، با بایاس b_o که آن هم در طول

¹ bias

² element-wise

یادگیری، به دست آمده است جمع گردد و نهایتاً از یک تابع فعالسازی sigmoid عبور کند. رابطه‌ی (۵-۳) بیانگر این عملیات است.

$$o_t = \sigma(w_o[[h_{t-1}, x_t] + b_o) \quad (۵-۳)$$

حال به سراغ محاسبه‌ی خروجی می‌رویم. پس از آنکه سلول حافظه در زمان t را از تابع فعالسازی \tanh عبور دادیم، با ضرب عنصر به عنصر دروازه‌ی 0 در آن، خروجی در این گام یعنی h_t به دست می‌آید. فرمول به دست آوردن خروجی در رابطه‌ی (۶-۳) آمده است.

$$h_t = o_t * \tanh(C_t) \quad (۶-۳)$$

شبکه‌های LSTM گرچه مشکل محو گرادیان را تا حد خوبی حل می‌کنند و اطلاعات گذشته را مدت زمان بیشتری در حافظه‌ی خود نگه می‌دارند، اما همچنان مشکلاتی هست که توسط این شبکه‌ها نیز حل نمی‌شود. مثلاً با به کار بردن این معماری نیز نمی‌توان پردازش داده‌ها را به صورت موازی پیش برد و پردازش حتماً باید به صورت ترتیبی انجام شود. بعلاوه دیدیم که معماری LSTM نسبت به RNN، پیچیده‌تر است. این پیچیدگی بدان معناست که تعداد پارامترهای بیشتری نسبت به RNN دارد و آموزش این نوع از مدل‌ها، حتی از شبکه‌های RNN نیز کندتر انجام می‌شود. هم‌چنین گرچه مدت‌زمانی که LSTM می‌تواند اطلاعات گذشته را حفظ کند نسبت به RNN پیشرفت داشته است، با اینهمه نمی‌توان گفت LSTM از پس حفظ وابستگی‌های بلندمدت نیز به خوبی برمی‌آید.

چالش‌هایی که برشمردیم، باعث شدند تا مدل‌های دیگری پدید آیند که مبتنی بر مبدل‌ها هستند. با این مدل‌ها در فصل آینده آشنا خواهیم شد.

فصل چهارم

روش‌های مبتنی بر مبدل

چندین سال است که از مفهوم یادگیری انتقالی در حوزه‌ی بینایی ماشین استفاده می‌شود. یادگیری انتقالی به این معناست که شبکه‌ای عصبی را روی دادگانی حجیم، به صورت بدون نظارت برای انجام کار مشخصی آموزش می‌دهیم و سپس از این مدل به عنوان لایه‌های پایه‌ای برای تنظیم دقیق روی داده‌های نظارت‌شده‌ی خود، استفاده می‌کنیم. مبدل، در واقع تلاشی برای پیاده‌سازی این مفهوم در زمینه‌ی پردازش زبان طبیعی است. در این فصل ابتدا به توضیح این مدل پرداخته و پس از آن دو مدل مبتنی بر مبدل شامل BERT و ParsBERT را معرفی خواهیم کرد.

۱-۴ مبدل

مدل مبدل، با تکیه بر مکانیسم توجه^۱، پردازش موازی داده‌های ترتیبی را مانند متن، فراهم آورده است. این موضوع افزایش سرعت چشمگیری در فرآیند آموزش این مدل‌ها را باعث شده است. غیر از این، امکان استفاده از مفهوم یادگیری انتقالی نیز با استفاده از مدل‌های مبتنی بر مبدل، میسر گردیده است. [۱۴] در ادامه به توضیح این مدل می‌پردازیم.

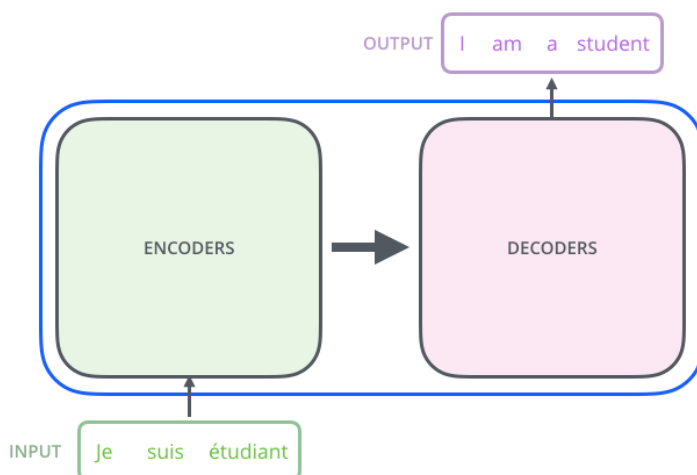
¹ attention mechanism

۱-۱-۴ معماری مبدل

در این بخش با معماری مدل بیشتر آشنا می‌شویم.

۱-۱-۴-۱ معماری رمزگذار-رمزگشا

نسخه‌ی اولیه‌ی این مدل بر روی کار ترجمه‌ی ماشینی انجام گرفته است. لذا در این بخش، با در نظر داشتن این کار، به تشریح معماری مدل می‌پردازیم. یک مدل مناسب برای ترجمه‌ی ماشینی، معماری رمزگذار^۱-رمزگشا^۲ است. این معماری برای کار ترجمه‌ی ماشینی در تصویر (۱-۴) آمده است. در این مثال، زبان مبدأ، فرانسه، زبان مقصد، انگلیسی و جمله‌ی مورد ترجمه «Je suis étudiante» به معنای «من دانشجو هستم» در نظر گرفته شده است.



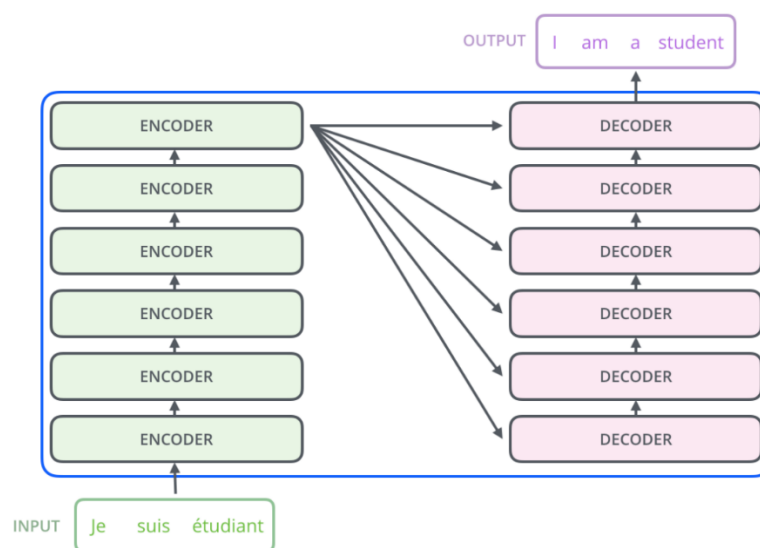
شکل (۱-۴) معماری رمزگذار-رمزگشا در ترجمه‌ی ماشینی [۶]

بخش رمزگذار، از شش لایه رمزگذار که روی هم پشته شده‌اند، تشکیل شده است. در مورد بخش رمزگشا هم همین شرایط برقرار است. در هر بخش، خروجی لایه‌ی قبلی، ورودی لایه‌ی بعدیست. اتصال بین رمزگذارها و رمزگشاها به شکل کلی، در تصویر (۲-۴) آمده است.

^۱ encoder

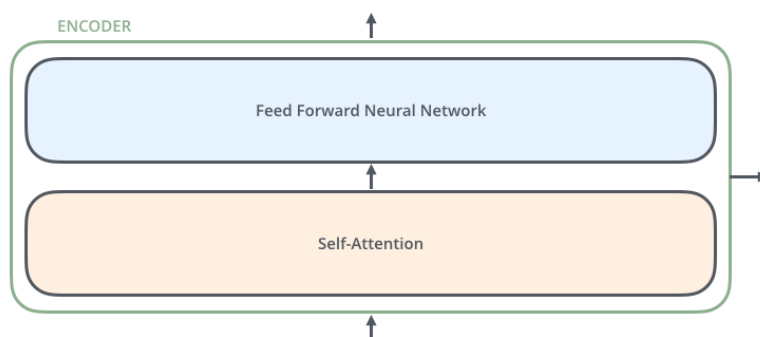
^۲ decoder

^۳ stack



شکل (۲-۴) اتصال دو بخش رمزگذار و رمزگشا در مبدل [۶]

هر رمزگذار از دو زیرلایه‌ی خودتوجهی^۱ و شبکه‌ی عصبی پیش‌خور^۲ تشکیل شده است. همچنین از یک اتصال باقی‌مانده در اطراف هر یک از زیرلایه‌ها استفاده می‌شود و یک لایه نرمالسازی پس از آن قرار می‌گیرد. شکل (۳-۴) دو زیرلایه اصلی را بدون اتصال باقی‌مانده و لایه نرمالسازی، برای هر رمزگذار نشان می‌دهد.



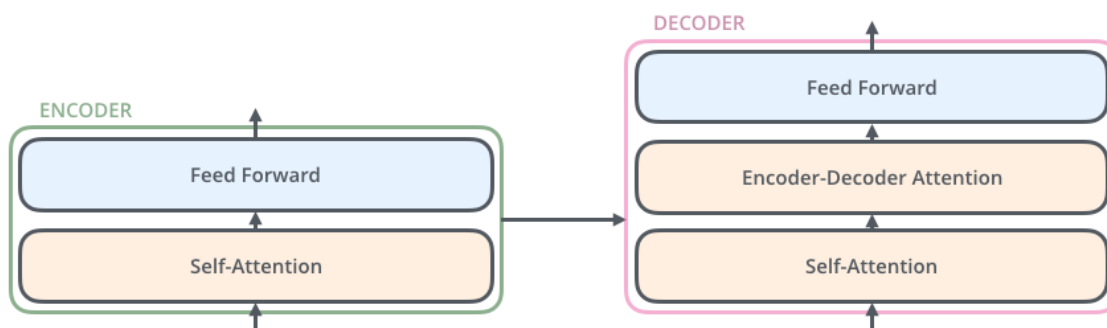
شکل (۳-۴) زیرلایه‌های اصلی رمزگذار [۶]

¹ self attention

² Feed Forward Neural Network (FFNN)

³ residual connection

هر رمزگشا، علاوه بر زیرلایه‌های ذکر شده برای رمزگذار، یک زیرلایه‌ی دیگر نیز به نام توجه رمزگذار-رمزگشا^۱ دارد. این لایه مکانیسم توجه چندسره^۲ را روی خروجی رمزگذار اعمال می‌کند. مشابه رمزگذار، در بخش رمزگشا نیز یک اتصال باقی‌مانده در اطراف هر زیرلایه وجود دارد و به دنبال آن یک لایه نرمالسازی انجام می‌شود. شکل (۴-۴) دو بخش رمزگشا و رمزگذار را با زیرلایه‌های اصلی آن‌ها نشان می‌دهد.



شکل (۴-۴) زیرلایه‌های اصلی رمزگذار و رمزگشا [۶]

گفتیم که ورودی هر یک از رمزگذارها، خروجی رمزگذار قبلی است. به این ترتیب کلمات ورودی پس از آنکه توسط یک الگوریتم تعبیه‌سازی^۳ به بردارهایی نظیر شدند، وارد رمزگذار اول می‌شوند. هر یک از کلمات ورودی به شکل برداری با اندازه‌ی ۵۱۲ نمایش داده می‌شوند و تعداد این کلمات بسته به طول متن می‌تواند متفاوت باشد. به شکل معمول این تعداد را برابر با طول بزرگترین جمله‌ی ورودی در نظر می‌گیرند. در ادامه برای سادگی کار، اندازه‌ی هر بردار را روی شکل چهار فرض می‌کنیم. شکل (۵-۴) نشان‌دهنده‌ی همین موضوع است. x_i به بردار نظیر Je کلمه اطلاق می‌شود.



شکل (۵-۴) نحوه‌ی نمایش هر کلمه به شکل بردار با اندازه‌ی فرضی [۶]

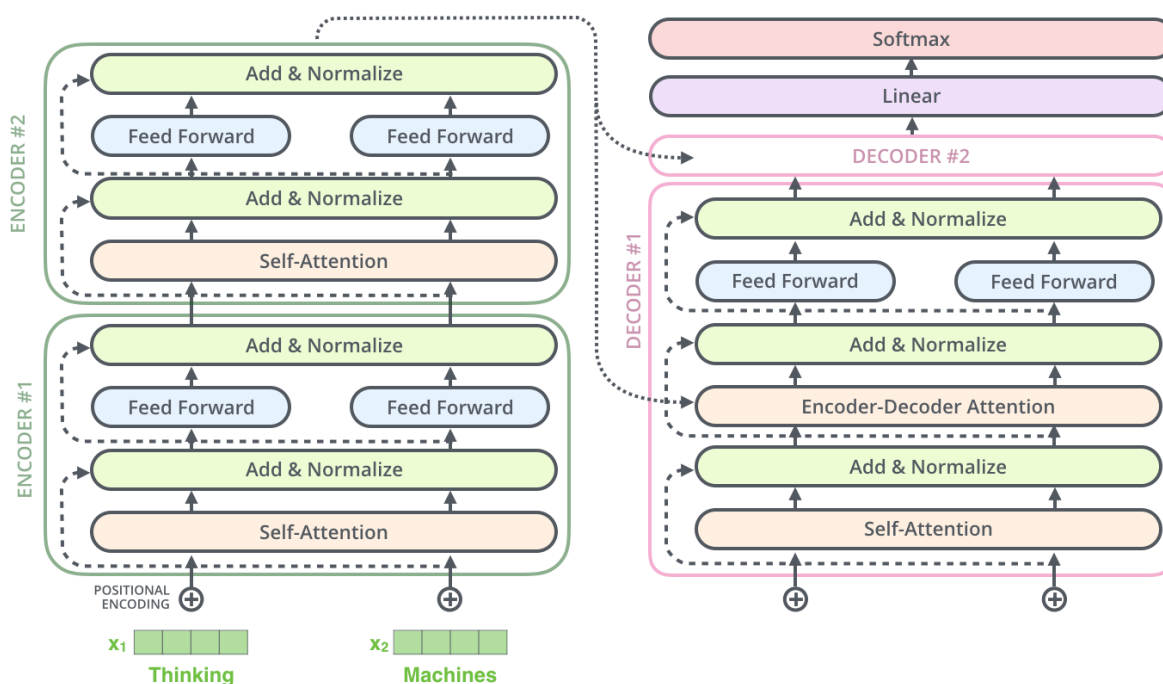
خروجی آخرین رمزگشا در هر مرحله یک بردار خواهد بود اما آنچه ما می‌خواهیم به دست آوریم، یک کلمه در زبان مقصد است. لذا برای آنکه از این بردارها، کلمات را استخراج کنیم، ابتدا به یک لایه شبکه عصبی خطی نیاز داریم که بردار ما را به برداری با اندازه‌ی بزرگتر نظیر کند. برای مثال اگر تعداد لغاتی که از پیش آموخته شده است، هزار باشد،

^۱ encoder-decoder attention

^۲ multi-head attention

^۳ embedding

اندازه‌ی برداری که پس از اعمال لایه‌ی خطی خواهیم داشت، باید هزار باشد. در این صورت هر درایه از بردار را می‌توان متناظر با یک کلمه از کل لغات در نظر گرفت. نهایتاً با اعمال softmax این امتیازات، به احتمال تبدیل می‌شوند و بالاترین احتمال، به عنوان کلمه‌ی خروجی ما در نظر گرفته می‌شود. در نتیجه معماری کلی این مدل با فرض آنکه تعداد رمزگشاها و رمزگذارها، ۲ باشد، در شکل (۴-۶) آمده است:



شکل (۴-۶) معماری کلی مدل با فرض وجود دو رمزگذار و رمزگشا [۶]

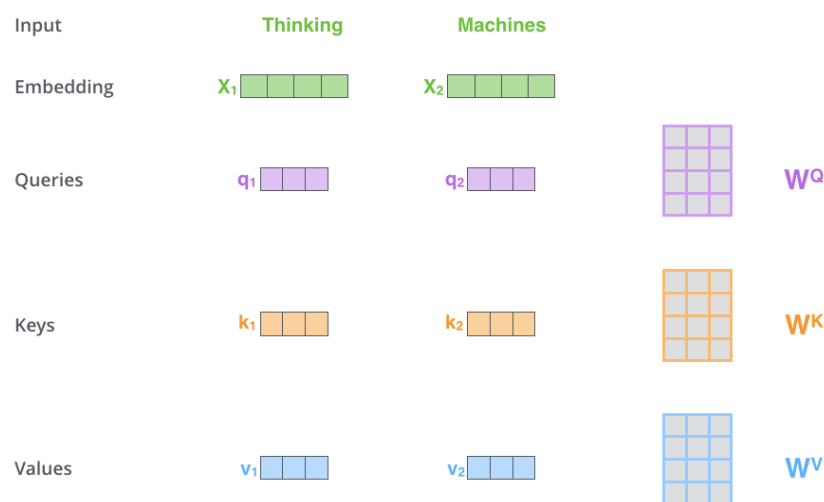
۴-۱-۱-۲ مکانیسم توجه

همانطور که گفته شد، با استفاده از مفهوم توجه، مدل مبدل به نتایج چشمگیری در زمینه‌ی پردازش زبان طبیعی دست یافت. در این بخش می‌خواهیم در مورد این مکانیسم و چگونگی اعمال آن در مبدل‌ها صحبت کنیم.

به شکل ساده مکانیسم توجه بیان می‌کند در تولید خروجی، به چه بخش‌هایی از ورودی باید بیشتر توجه کنیم. این مکانیسم باعث می‌شود مدل‌ها در کارهای مربوط به پردازش زبان طبیعی مانند سیستم‌های پرسش و پاسخ، بهتر عمل کنند. در مورد این مثال، توجه به این شکل معنی می‌شود که در پاسخ به یک پرسش از روی متن، روی چه بخش‌هایی از متن باید تمرکز بیشتری کرد. در مورد کار ترجمه‌ی ماشینی هم می‌توان گفت مکانیسم توجه به مدل کمک می‌کند تا بتواند تصمیم بگیرد در ترجمه‌ی هر کلمه، به چه کلماتی از ورودی باید بیشتر توجه کند.

در این مکانیسم ابتدا برای هر کلمه‌ی ورودی که با یک بردار نمایش داده می‌شود (x_i) ، سه بردار با نام‌های Key، Query و Value ساخته می‌شود. بردار Query مربوط به کلمه‌ی نام را با q_i ، بردار Key مربوط به کلمه‌ی نام را با k_i و بردار Value مربوط به کلمه‌ی نام را با v_i نشان می‌دهیم. سائز این سه بردار معمولاً از بردار اصلی کوچک‌تر است. در مدل اصلی اندازه‌ی بردار هر کلمه ۵۱۲ و اندازه‌ی مربوط به هر یک بردارهای ساخته‌شده برای آن، برابر با ۶۴ می‌باشد. این بردارها از ضرب بردار اصلی در سه ماتریس W^Q ، W^K و W^V که به همراه مدل، آموزش دیده‌اند، به دست می‌آید.

با در نظر گرفتن عبارت "Thinking Machines" به عنوان دنباله‌ی ورودی، سه بردار ذکر شده برای هر یک از کلمات Thinking و Machines ساخته می‌شود. در شکل (۷-۴) ابعاد این بردارها و ماتریس‌های W^Q ، W^K و W^V به صورت نمادین آمده است.



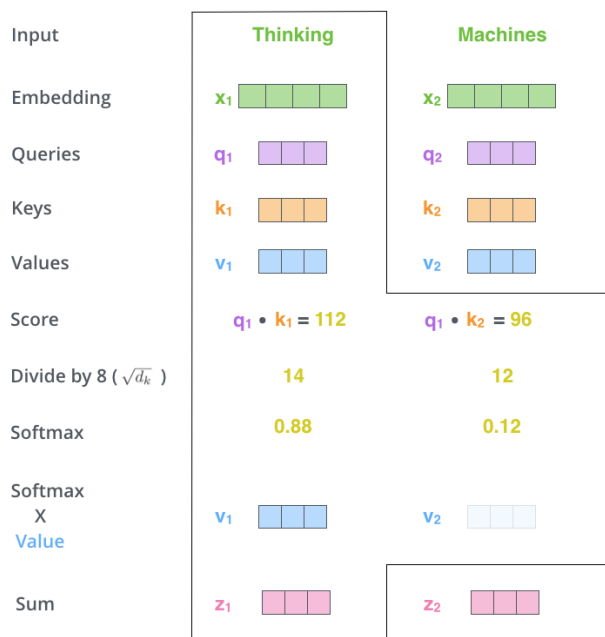
شکل (۷-۴) ابعاد نمادین ماتریس‌ها و بردارهای ساخته‌شده [۶]

در قدم بعدی امتیاز هر یک از کلمات در قیاس با کلمه‌ای که در حال پردازش آن هستیم محاسبه می‌شود تا بدانیم کدام بخش از ورودی شایسته‌ی توجه بیشتری است. در این گام اگر در حال پردازش اولین کلمه هستیم، ضرب داخلی بردار Query متناظر با آن را در تک‌تک بردارهای Key متناظر با تمام کلمات ورودی محاسبه می‌کنیم.

پس از آن هر یک از این امتیازات بر هشت تقسیم می‌شوند. عدد هشت از این رو انتخاب می‌شود که جذر اندازه‌ی بردارهای ساخته‌شده، یعنی ۶۴ می‌باشد. این کار برای پایدارسازی گرادیان‌ها صورت می‌گیرد. در نهایت با انجام عملیات softmax روی این اعداد، امتیازات نرمال شده و مجموع آن‌ها به یک می‌رسد. گام بعدی آنست که اعداد حاصل از عملیات softmax را در بردار Value نظیر هر کلمه‌ی ورودی ضرب کنیم. نهایتاً بردارهای به دست آمده را با یکدیگر جمع می‌کنیم تا خروجی این زیرلایه برای کلمه‌ی اول به دست آید. با این کار در واقع ما یک جمع وزن‌دار ترتیب داده‌ایم که کلماتی که امتیاز

بیشتری داشته‌اند، تاثیر بیشتری نیز داشته باشند و به لایه‌ی بالاتر انتقال پیدا کنند. بردار ساخته شده در این مرحله برای کلمه‌ی i ام را z_i می‌نامیم.

شکل (۴-۸) برای درک بهتر این عملیات و با اعداد فرضی آورده شده است.

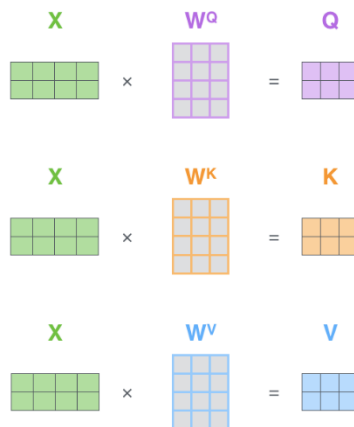


شکل (۴-۸) عملیات محاسبه‌ی بردار خروجی برای هر کلمه [۶]

این کار برای تمام کلمات ورودی به شکل گفته‌شده، انجام می‌گیرد. گرچه برای ساده شدن انتقال مطلب، در اینجا محاسبات را کلمه به کلمه در نظر گرفتیم، اما در واقع، این عملیات به صورت ماتریسی صورت می‌گیرد تا سرعت انجام محاسبات بیشتر شود. در این صورت ما به جای آنکه برای کلمه‌ی i ام بردارهای x_i ، q_i و k_i را داشته باشیم، برای تمام کلمات ورودی، ماتریس X ، Q ، K و V را خواهیم داشت. ابعاد این ماتریس‌ها با در نظر گرفتن فرضیات قبلی در شکل (۴-۹) آمده است. بنابراین نهایتاً خروجی تمام کلمات در ماتریسی به نام Z ذخیره خواهد شد.

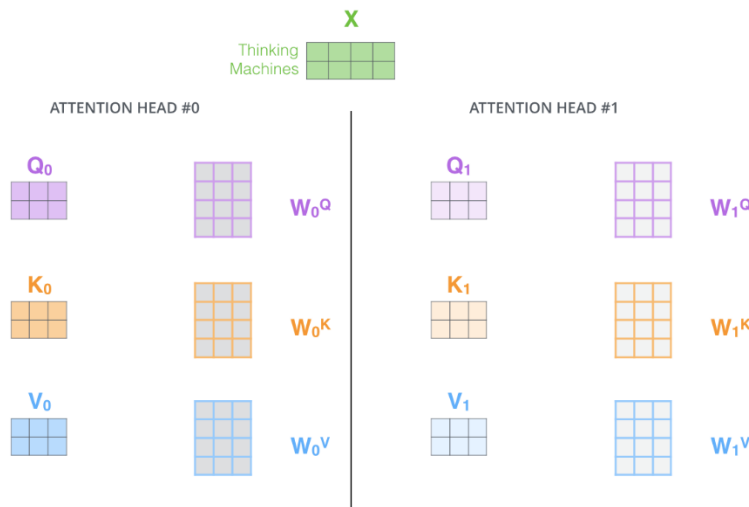
۴-۱-۱-۳ مکانیسم توجه چندسر

دیدیم که در محاسبه‌ی امتیازات برای کلمه‌ی در حال پردازش، خود آن کلمه را نیز دخیل دانستیم و امتیاز آن را نیز محاسبه کردیم. گاهی پیش می‌آید که امتیازات باقی کلمات، مغلوب کلمه‌ی اصلی واقع شوند و مکانیسم توجه به خوبی عمل نکند. برای رفع این مشکل، رویکردی با نام چند سر در این مدل معرفی می‌شود.



شکل (۴-۹) فرم ماتریسی محاسبات [۶]

این مکانیسم به مدل این امکان را می‌دهد که بتواند روی بخش‌های مختلف ورودی تمرکز کند. در این مکانیسم در واقع ما چندین مجموعه Query، Key و Value خواهیم داشت که هم‌چون قبل، از ضرب ماتریس کلمات در ماتریس‌های آموزش دیده‌ی W^Q ، W^K و W^V به دست می‌آیند. برای هر سر از توجه، این ماتریس‌ها متفاوتند. شکل (۴-۱۰) این تفاوت را نمایش می‌دهد. در مدل مبدا اصلی، تعداد این سرها، هشت است.



شکل (۴-۱۰) تفاوت ماتریس‌ها در سرهای مختلف توجه [۶]

در این صورت، با محاسبه‌ی امتیازات به گونه‌ای که گفته شد، هشت ماتریس Z متفاوت به دست می‌آید. این هشت ماتریس را در کنار هم الحاق می‌کنیم تا ماتریس بزرگتری به دست آید. ماتریس به دست آمده را در ماتریس وزن W^0 که همراه با مدل، آموزش دیده است، ضرب می‌کنیم تا خروجی نهایی زیرلایه‌ی خودتوجهی به دست آید.

این عملیات در شکل (۱۱-۴) آمده است.

1) Concatenate all the attention heads

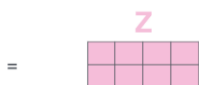


2) Multiply with a weight matrix W^O that was trained jointly with the model

\times



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



شکل (۱۱-۴) عملیات تولید خروجی نهایی زیرلایه‌ی خودتوجهی [۶]

۴-۱-۱-۳ کاربرد مکانیسم توجه در بخش‌های مختلف مبدل

در زیرلایه‌ی توجه رمزگذار-رمزگشا در سمت رمزگشا، ماتریس Q از لایه‌ی قبلی رمزگشا تامین می‌شود. ماتریس‌های K و V از خروجی آخرین رمزگذار می‌آیند. این مورد به رمزگشا اجازه می‌دهد که نسبت به تمام کلمات دنباله‌ی ورودی آگاهی داشته باشد و بداند در هر گام، کدام یک از آن‌ها به توجه بیشتری نیاز دارند.

در زیرلایه‌ی خودتوجهی در سمت رمزگذار، هر سه ماتریس Q ، K و V از خروجی لایه قبلی در رمزگذار می‌آیند. به این دلیل به آن، مکانیسم خودتوجهی گفته می‌شود. در واقع هر موقعیت^۱ در رمزگذار می‌تواند به تمام موقعیت‌های لایه قبلی رمزگذار توجه کند.

در زیرلایه‌ی خودتوجهی در سمت رمزگشا، به طور مشابه، به هر موقعیت در رمزگشا اجازه داده می‌شود که تمام موقعیت‌های رمزگشا از ابتدا تا آن موقعیت، توجه کند.

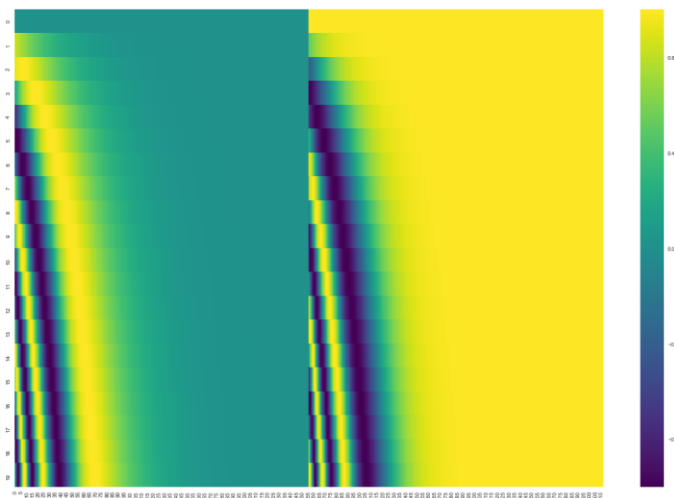
^۱ position

۴-۱-۲ رمز گذاری موقعیتی

می دانیم که در داده های زبانی، ترتیب اهمیت بالایی دارد؛ بعنوان مثال دو عبارت «کار برای زندگی» و «زندگی برای کار»، با آنکه متشکل از تعدادی کلمه ی کاملاً مشترک هستند، به دلیل تفاوت در موقعیت کلمات، معانی کاملاً متفاوتی را منتقل می کنند.

تا اینجا در مدل مبدل، دیدیم که کلمات به شکل موازی پردازش می شوند. سوالی که پیش می آید اینست که فاکتور تاثیر گذاری مثل ترتیب چگونه در این مدل ها لحاظ می گردد؟ باید گفت این مدل از رمز گذاری موقعیتی^۱ استفاده می کند. به این معنا که جایگاه هر کلمه در مرحله ی تعبیه سازی به همراه خود کلمه در بردار نهایی آن واژه، تعبیه می شود و هر بردار علاوه بر محتوای کلمه، اطلاعاتی از جایگاه کلمه در جمله را نیز دارد. بدین ترتیب بردار دیگری که دارای الگویی خاص است و مربوط به جایگاه کلمات می شود، به بردار مربوط به کلمه اضافه می شود. این بردار از فرمولی که در مقاله ی اصلی آمده است، محاسبه می شود که به علت پیچیدگی آن، از آوردن این رابطه در گزارش خودداری کرده ایم.

شکل (۴-۱۲) مربوط به این الگو برای دنباله ای از بیست کلمه با سایز ۵۱۲ می باشد. همانطور که مشخص است، نیمه ی سمت راست تصویر و نیمه ی سمت چپ آن به شکل مشهودی، متفاوتند. این بدان دلیل است که یک نیمه از رابطه ی کسینوسی و نیمه ی دیگر از رابطه ی سینوسی به دست آمده و نهایتاً کنار یکدیگر الحاق شده اند. برای مثال اگر بخواهیم بردار نهایی کلمه ی اول را برای ورود به رمز گذار اول به دست آوریم، آن را با برداری که ۲۵۶ عنصر اول آن صفر و ۲۵۶ عنصر بعدی آن یک است جمع می کنیم.



شکل (۴-۱۲) بردارهای رمز گذاری موقعیتی برای بیست کلمه با اندازه ی ۵۱۲ [۶]

^۱ positional encoding

۲-۴ مدل BERT

مدل BERT^۱ یک مدل زبانیست که در سال ۲۰۱۸ و توسط گوگل معرفی شد. در زمان انتشار این مدل، نتایج مربوط به یازده کار حوزه‌ی پردازش زبان طبیعی، از جمله پرسش و پاسخ با داده‌های SQuAD، بهبود قابل قبولی پیدا کردند. [۳]

این مدل یک مدل دوطرفه‌ی چندلایه از رمزگذارهای مبدل‌هاست که در بخش قبل با آن آشنا شدیم. پیش از انتشار مدل BERT، مدل‌هایی مانند ELMo^۲ و OpenAI GPT کاربرد داشتند که هر یک، با مشکلاتی مواجه بودند.

از مدل زبانی ELMo می‌توان برای تعبیه‌سازی کلمات به روشی جدید استفاده کرد که برخلاف الگوریتم‌های تعبیه‌سازی پیشین همچون GloVe [۹] یا word2vec [۷]، کلمه را با توجه به کاربرد آن در جمله تعبیه می‌کند. برای مثال برای کلمه‌ی stick چندین تعبیه‌سازی مختلف وجود دارد که به فراخور بافت متنی که این کلمه در آن ظاهر می‌شود، انتخاب می‌شود که کدام تعبیه‌سازی مناسب است. این مدل مبتنی بر LSTM ساخته شده است و ایراداتی که برای این نوع شبکه‌ها برشمردیم، متوجه این مدل نیز هست. علاوه بر این، با وجود آنکه LSTM به شکل دوطرفه در معماری آن به کار رفته است، اما به این دلیل که صرفاً نتایج حاصل از مسیر رفت و برگشت را به یکدیگر الحاق می‌شوند، نمی‌توان گفت مدل از هر دو جهت درک خوبی به دست می‌آورد. [۱۰]

مدل OpenAI GPT با اینکه مبتنی بر مبدل‌ها ساخته شده است و مشکلات شبکه‌های عصبی بازگشتی را ندارد، اما به دلیل آنکه تنها در یک جهت عمل می‌کند، نمی‌تواند از هر دو سمت کلمات، آگاهی کسب کند. این مدل از چند لایه رمزگشای مبدل‌ها که به صورت پشت‌پشته قرار گرفته‌اند، تشکیل شده است. [۱۱]

مدل BERT، هیچ یک از مشکلات دو مدل یادشده را ندارد و همانطور که گفته شد، هم مبتنی بر مبدل‌هاست و هم از یک معماری دوطرفه برخوردار است. در ادامه با معماری این مدل آشنا می‌شویم.

۱-۲-۴ معماری BERT

این مدل در دو نسخه با اندازه‌های متفاوت ارائه شده است. برای هر یک از آن‌ها، تعداد لایه‌های رمزگذار، اندازه‌ی حالت پنهان، تعداد سرهای توجه و تعداد کل پارامترها در جدول زیر آمده است:

تعداد کل پارامترها	تعداد سرهای توجه	اندازه‌ی حالت پنهان	تعداد لایه‌های رمزگذار	
110 M	12	768	12	BERT BASE
340 M	16	1024	24	BERT LARGE

جدول (۱-۴) تعداد اجزای مربوط به معماری دو نسخه‌ی BERT

¹ Bidirectional Encoder Representations from Transformers

² Embeddings From Language Models

۲-۲-۴ نمایش ورودی در مدل BERT

مدل BERT به گونه‌ای طراحی شده است که ورودی آن هم می‌تواند یک جمله‌ی واحد باشد و هم دو جمله‌ی مجزا از هم. باید توجه داشت منظور از جمله در اینجا، لزوماً یک جمله‌ی زبانی نیست؛ دنباله‌ی به هم پیوسته‌ای از کلمات را در این مدل می‌توان یک جمله در نظر گرفت. ترجمه‌ی ماشینی می‌تواند مثالی از یک کار پایین‌دستی^۱ باشد که ورودی آن تنها یک جمله است. در حالیکه سیستم‌های پرسش و پاسخ، برای استخراج جواب، به دو جمله‌ی پرسش و متن به عنوان ورودی نیاز دارند.

در مدل BERT، برای تعبیه‌سازی کلمات، از الگوریتم WordPiece استفاده می‌شود. این الگوریتم، تعداد سی هزار نشانه^۲ را در خود دارد. در این نوع از تعبیه‌سازی، کلمات ممکن است خود به قسمت‌های کوچک‌تر تنظیم شوند. برای مثال اگر بخواهیم واژه‌ی snowing را تعبیه کنیم، توسط این الگوریتم کلمه به دو بخش snow و ing^۳ شکسته می‌شود و هر یک به شکل مستقل تعبیه می‌گردند.

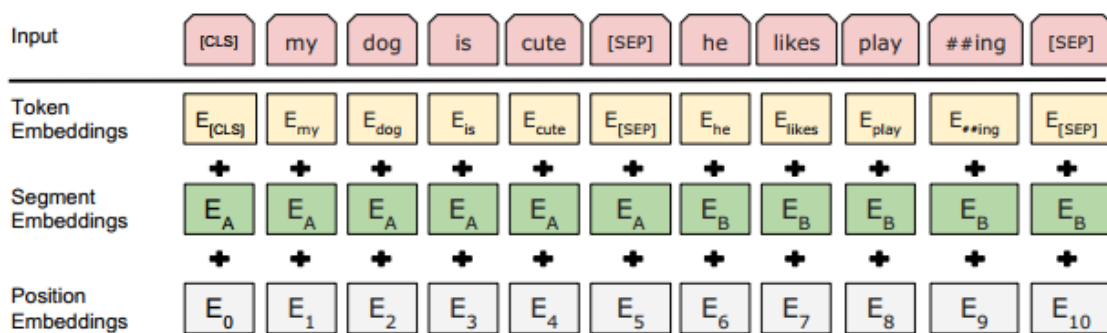
برای آماده کردن دنباله‌ی ورودی برای ورود به مدل، ابتدای کار یک توکن خاص به نام [CLS] که مخفف Classification است به ابتدای رشته‌ی ورودی اضافه می‌شود. در انتهای جمله نیز توکن خاص [SEP] ظاهر می‌شود که مخفف Seperate است؛ این توکن برای جدا کردن جملات نامتوالی به کار می‌رود. اگر نمونه‌های ما تک جمله‌ای باشند، نشانه‌ی آخر آن‌ها و اگر دو جمله‌ای باشند، نشانه‌ی میان دو جمله و نشانه‌ی انتهایی آن‌ها، [SEP] خواهد بود. جملات، جز آنکه با نشانه‌ی از هم جدا می‌شوند، از طریق روش دیگری نیز این کار برایشان صورت می‌گیرد. به این صورت که بردار هر کلمه، با برداری که تعیین می‌کند این کلمه متعلق به جمله‌ی اول است یا جمله‌ی دوم، جمع می‌شود. به این کار تعبیه‌سازی بخش گفته می‌شود.

در نتیجه برای هر کلمه‌ی ورودی، بردار مربوط به کلمه در WordPiece، بردار تعبیه‌سازی موقعیتی کلمه که در بخش مبدل به آن پرداختیم و بردار تعبیه‌سازی بخش که مربوط به جدا کردن دو جمله از یکدیگر است، با یکدیگر جمع می‌شوند تا بردار نهایی ورودی هر کلمه را بسازند. شکل (۴-۱۳) برای نشان دادن این عملیات آمده است.

^۱ downstream task

^۲ token

^۳ segment embedding



شکل (۴-۱۳) عملیات آماده‌سازی نشانه‌های ورودی [۳]

۴-۲-۳ رویکردهای استفاده از مدل BERT

در سیستم‌هایی که از BERT استفاده می‌کنند، می‌توان از دو رویکرد استفاده کرد. این رویکردها عبارتند از: رویکرد مبتنی بر ویژگی و رویکرد تنظیم دقیق.

در رویکرد مبتنی بر ویژگی، خروجی ما کلمات تعبیه‌شده هستند اما تفاوت این تعبیه‌سازی‌ها با الگوریتم‌های معروفی مانند GloVe یا word2vec در اینست که در BERT، کلمات با در نظر داشتن معانی آن‌ها در بافت جمله درون‌سازی می‌شوند و کلمات در بافت‌های متنی متفاوت می‌توانند نمایش متفاوتی داشته باشند. این تعبیه‌ها را می‌توان در یک مدل از پیش ساخته‌شده بعنوان ورودی استفاده کرد.

در رویکرد تنظیم دقیق، مدلی که با داده‌گان بدون برچسب ساخته‌ایم را برای تنظیم دقیق یک کار نظارت‌شده‌ی پایین‌دستی به کار می‌بریم. مزیت این استراتژی آن است که تعداد کمی از پارامترها را باید از ابتدا یاد گرفت و اکثراً از پارامترهای از پیش آموزش دیده استفاده می‌شود. این مورد باعث می‌شود بتوان در هزینه و انرژی صرفه‌جویی چشمگیری کرد. به طوریکه در طول ۳-۴ اپیاک هم می‌توان مدل را برای کار مورد نظر آماده کرد.

با وجود اینکه از این مدل می‌توان با هر دو استراتژی یادشده استفاده کرد، اما باید گفت ارزش افزوده‌ای که BERT ایجاد کرده است، بیشتر در خدمت رویکرد تنظیم دقیق است. لذا در بخش بعد به طور دقیق‌تر به رویکرد تنظیم دقیق می‌پردازیم.

۴-۲-۴ استفاده از مدل BERT با رویکرد تنظیم دقیق

این رویکرد از دو مرحله‌ی اصلی ساخته شده است که عبارتند از: پیش‌آموزش و تنظیم دقیق.

¹ feature-based approach

² pre-training

در طول مرحله‌ی پیش‌آموزش، مدل بر روی داده‌های بدون برچسب، آموزش داده می‌شود. برای تنظیم دقیق، ابتدا مدل BERT با پارامترهای از پیش آموزش‌داده‌شده مقداردهی اولیه می‌شود و همه‌ی پارامترها با استفاده از داده‌های برچسب‌گذاری شده از وظایف پایین‌دستی تنظیم می‌شوند. هر کار پایین‌دستی گرچه با پارامترهای یکسان مقداردهی اولیه می‌شود، اما دارای مدل‌های تنظیم‌شده‌ی جداگانه‌ای است.

ویژگی متمایزکننده‌ی BERT معماری یکپارچه آن در مورد کارهای مختلف است. چندانکه بین معماری مدل از پیش‌آموزش دیده و مدل نهایی برای یک کار پایین‌دستی، حداقل تفاوت وجود دارد.

۴-۲-۴ مرحله‌ی پیش‌آموزش

همانطور که پیش‌تر گفته شد، مدل BERT یک مدل زبانی دوطرفه است. ضمناً می‌دانیم تعریف مدل زبانی به صورت ساده، پیش‌بینی کلمه‌ی بعدی در دنباله‌ای از داده‌ها، در نظر گرفته می‌شود. با در نظر داشتن این دو نکته، سوالی که مطرح می‌شود آنست که اگر بخواهیم هر کلمه‌ای، از دادگان پیش و پس از خود استفاده کند و دوطرفه باشد، آن کلمه به شکل غیرمستقیم خود را خواهد دید و به شکلی بدیهی پیش‌بینی درستی خواهد داشت که مورد استفاده‌ی ما نیست. پس چطور می‌توان معماری را به شکل دوطرفه طراحی کرد؟

راه حلی که طراحان مدل BERT برای این مشکل ارائه داده‌اند آنست که مدلی طراحی کنند که درصد مشخصی از کلمات ورودی را به صورت تصادفی پنهان کند و سپس به حدس آن کلمات پردازد. از این تکنیک با عنوان Masked Language Model یا به اختصار MLM یاد می‌شود و به این ترتیب آموزش دوطرفه ممکن می‌شود.

تعداد خوبی از کارهای پایین‌دستی مانند سیستم‌های پرسش و پاسخ و یا استنتاج زبان طبیعی، به درک رابطه‌ی میان دو جمله نیاز دارند اما به نظر می‌رسد MLM به تنهایی نمی‌تواند این درک را ممکن سازد. در نتیجه کار دیگری به نام Next Sentence Prediction یا پیش‌بینی جمله‌ی بعد که به اختصار به آن NSP گفته می‌شود، به عنوان دومین گام که در مرحله‌ی پیش‌آموزش انجام می‌شود، به میان آمده است تا این مشکل را حل کند.

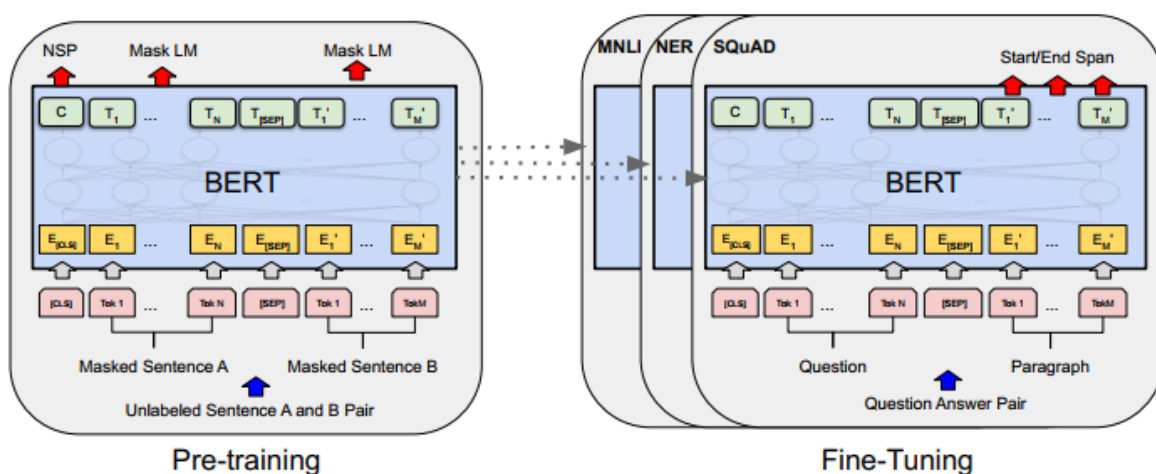
در این تکنیک، دو جمله به مدل داده می‌شود و مدل باید یاد بگیرد آیا جمله‌ی دوم، دنباله‌ی جمله‌ی اول بوده است یا خیر. به این صورت که در نیمی از دادگان پیش‌آموزش، دو جمله‌ای که مدل دریافت می‌کند دنباله‌ی هم هستند و در نیمی دیگر کاملاً بی‌ربط بوده و دنباله‌ی یکدیگر نیستند.

گفتنیست داده‌های مرحله‌ی پیش‌آموزش از ویکی‌پدیای انگلیسی و همچنین مجموعه‌داده‌ی BookCorpus به دست آمده است.

۲-۴-۲-۴ مرحله‌ی تنظیم دقیق

چه داده‌ی ورودی ما یک جمله‌ای باشد و چه دو جمله‌ای، مکانیسم توجه این امکان را فراهم می‌کند تا مدل بتواند به درک خوبی از مفهوم زبان برسد. در صورتی که ورودی ما از دو جمله تشکیل شده باشد، با استفاده از همین مکانیسم می‌توان به ارتباط میان آن دو نیز توجه کرد.

برخلاف مرحله‌ی پیش‌آموزش که به زمان زیادی برای آموزش نیاز داشتیم، در این مرحله، تنها با چند اپیاک می‌توان مدل را برای یک کار پایین‌دستی آماده کرد. شکل (۴-۱۴) این دو مرحله را در کنار هم نشان می‌دهد. همانطور که می‌بینیم در سمت چپ که مربوط به مرحله‌ی پیش‌آموزش است، دو کار مدل زبانی پنهان (MLM) و پیش‌بینی جمله‌ی بعدی (NSP) آمده‌اند. جملات A و B که درصدی از نشانه‌هایشان پنهان شده است، وارد فاز پیش‌آموزش می‌شوند و هر یک از کلماتشان، وارد بدنه‌ی اصلی مدل می‌گردند. پس از آنکه از ۱۲ یا ۲۴ لایه رمزگذار، بسته به اندازه‌ی مدل، عبور کردند، خروجی نهایی مربوط به هر کلمه حاصل می‌شود. دقت داریم که در ابتدای فرآیند، نشانه‌های خاص [CLS] و [SEP] اضافه شده‌اند. سپس تمام کلماتی که پنهان شده بودند حدس زده می‌شوند و خروجی C، اطلاعاتی درباره‌ی تمام ورودی را در اختیار ما می‌گذارد که با یک کار کلاس‌بندی دودویی، خروجی NSP از آن استخراج می‌شود که نشان می‌دهد آیا جمله‌ی B در ادامه‌ی جمله‌ی A بوده است یا خیر. در مرحله‌ی تنظیم دقیق، اگر کار را پرسش و پاسخ در نظر بگیریم، جمله‌ی اول پرسش و جمله‌ی دوم متن خواهد بود. به یاد داریم که منظور از جمله، جمله‌ی زبانی نیست که دارای تنها یک فعل باشد و هر متنی که قبل از نشانه‌ی [SEP] بیاید را می‌توان یک جمله در نظر گرفت. برای این کار پایین‌دستی پس از انجام فرآیند تنظیم دقیق، بازه‌ی شروع و پایان جواب به دست می‌آید.



شکل (۴-۱۴) مراحل پیش‌آموزش و تنظیم عمیق مدل BERT [۳]

۳-۴ مدل ParsBERT

گفتیم که مدل BERT یک مدل زبانی است که برای زبان انگلیسی طراحی گردیده است. اما اگر بخواهیم از این مدل برای کارهایی که با دادگان فارسی سروکار دارند، استفاده کنیم، چندان موثر نخواهد بود؛ چرا که به شکل طبیعی ساختار زبان فارسی و انگلیسی با یکدیگر تفاوت‌های اساسی دارند. در این بخش به معرفی مدل زبانی ParsBERT خواهیم پرداخت که با کمک مدل BERT انگلیسی، برای زبان فارسی طراحی شده است. [۴]

ParsBERT یک مدل زبانی مختص زبان فارسی است که در سال ۲۰۲۱ ارائه شد. معماری این مدل با الهام از معماری BERT انگلیسی و مشابه آن می‌باشد. به همین دلیل، در این فصل تنها به فرآیند آماده‌سازی داده‌های فارسی برای مدل، بسنده کرده‌ایم و باقی، شبیه به مدل BERT عمل می‌کند.

۱-۳-۴ جمع‌آوری داده

منابعی که داده‌های مرحله‌ی پیش‌آموزش مدل ParsBERT از آن‌ها استخراج شده است، در جدول (۲-۴) فهرست گردیده‌اند.

منبع	تعداد جمله
Persian Wikipedia	1878008
BigBang Page	3017
Chetor	166312
EliGasht	214328
DigiKala	177357
Ted Talks	46833
Books	25335
Miras-text	35758281

جدول (۲-۴) منابع دادگان مرحله‌ی پیش‌آموزش برای مدل ParsBERT

۲-۳-۴ پیش‌پردازش داده

پس از جمع‌آوری مجموعه داده‌ی پیش‌آموزش، چند مرحله پیش‌پردازش، از جمله پاکسازی^۱، جایگزینی و نرمال‌سازی، برای تبدیل مجموعه داده‌ها به یک قالب مناسب ضروری است. ابتدا باید تمام کاراکترهای ناخواسته و بی‌اهمیت موجود در مجموعه را پاک کنیم. این کار از طریق یک سری فرآیند انجام می‌شود که عبارتند از:

- (۱) اصلاح تمام کاراکترهای Unicode.
- (۲) حذف تمام اطلاعات بی‌اهمیت به عنوان مثال URL ها.
- (۳) حذف تمام برچسب‌های غیر ضروری مانند تگ‌های HTML، CSS و JS.

^۱ sanitizing

۴) نرماسازی متن فارسی با استفاده از کتابخانه‌ی HAZM که یک کتابخانه‌ی محبوب برای پیش پردازش متن فارسی است. نرمالسازی متن، اصلاح حروف فارسی، حذف فاصله‌های اضافی، تبدیل اعداد انگلیسی و عربی به فارسی، اصلاح علائم نگارشی و اصلاح نیم‌فاصله را شامل می‌شود.

۵) حذف کاراکترهای اسکی ناخواسته در زبان فارسی.

۶) جایگزینی کاراکترهای عربی با معادل آن‌ها در فارسی.

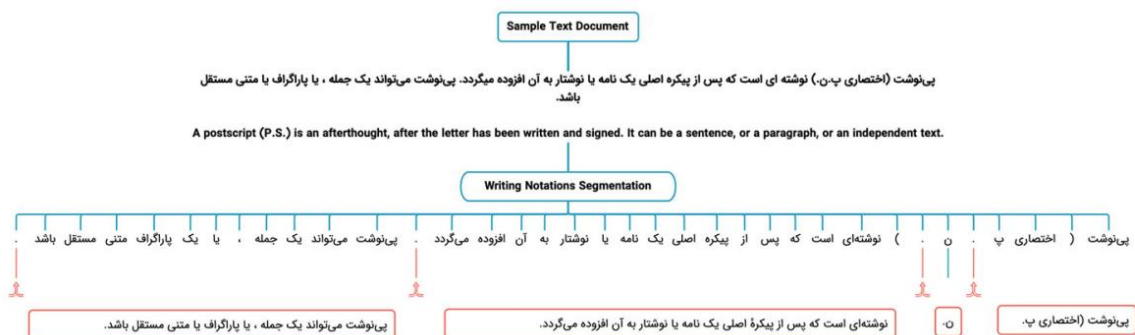
۷) حذف یا جایگزینی کاراکترهای Unicode بی‌معنی.

۴-۳-۳ تقسیم اسناد به جملات صحیح فارسی

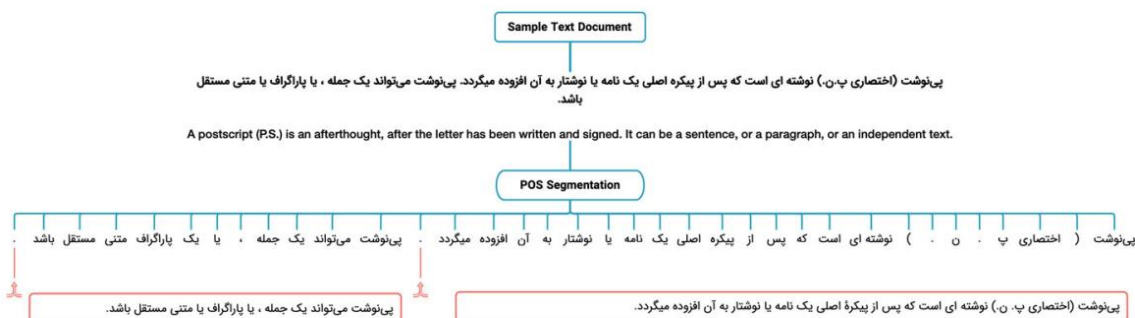
پس از اینکه مجموعه پیش‌پردازش شد، باید به جملات واقعی هر سند، تبدیل شود. یک جمله‌ی صحیح در فارسی بر اساس نمادهای نقطه، علامت تعجب، علامت سوال و دو نقطه تشخیص داده می‌شود. با این حال، تقسیم محتوا صرفاً بر اساس این نمادها، باعث بروز مشکلاتی می‌شود. شکل (۴-۱۵)، نمونه‌ای از چنین مسئله‌ای را نشان داده شده است. مشاهده می‌شود که نتیجه، شامل جملات کوتاه بی‌معنی است. زیرا در فارسی اختصاراتی وجود دارد که با علامت نقطه (.) از هم جدا شده‌اند. به عنوان یک جایگزین، استفاده از اطلاعات اجزای کلام^۱ می‌تواند راه حل مناسبی برای رسیدگی به این خطاها و تولید خروجی‌های دلخواه باشد. در زبان فارسی تمام جمله‌ها به فعل ختم می‌شوند. بنابراین، هر گاه قسمتی با نشانه‌ای که به عنوان «فعل» تگ شده است، پایان یابد و یکی از نمادها به دنبال آن قرار گیرد، به عنوان یک بخش معنی‌دار صحیح علامت‌گذاری شده و به عنوان یک جمله از بقیه قسمت‌ها جدا می‌شود.

پس از آنکه مجموعه داده‌ی پیش‌آموزش آماده شد، نوبت به دو مرحله‌ی پیش‌آموزش و تنظیم دقیق می‌رسد. این دو مرحله شبیه آنچه که در BERT دیدیم، اتفاق می‌افتد؛ بنابراین از تکرار آن‌ها پرهیز می‌کنیم. گفتنیست در مقاله‌ی اصلی، نتایج تنظیم دقیق این مدل برای سه کار تحلیل احساس، طبقه‌بندی متون و تشخیص موجودیت‌های نام‌دار آمده است. در فصل آینده در مورد تنظیم دقیق این مدل برای پرسش و پاسخ صحبت خواهیم کرد

^۱ Part Of Speech (POS)



(a)



(b)

شکل (۴-۱۵) مثالی از تقسیم یک متن به جملات با استفاده از دو شیوه (a): توجه به نشانه‌ها (b): توجه به اجزای کلام [۴]

فصل پنجم

حل مسئله

در این فصل به تنظیم دقیق مدل ParsBERT برای کار پرسش و پاسخ فارسی روی دادگان PersianQA، پیاده‌سازی مدل با استفاده از کتابخانه‌های اکوسیستم Hugging Face و نتایج نهایی مدل خواهیم پرداخت.

۵-۱ تنظیم دقیق ParsBERT برای کار پرسش و پاسخ

مدل BERT برای کار پرسش و پاسخ روی مجموعه داده‌ی SQuAD در هر دو نسخه، تنظیم دقیق شده است. به روشی مشابه می‌توان ParsBERT را روی PersianQA تنظیم دقیق کرد. هدف مدل ساخته شده با استفاده از هر یک از این دو مدل، آنست که بتواند بازه‌ی پاسخ را از روی متن داده شده پیدا کند.

در طول روند تنظیم دقیق، دو بردار شروع و پایان جواب معرفی می‌شوند که همراه مدل آموزش می‌بینند. اندازه‌ی هر یک از این بردارها به اندازه‌ی لایه‌ی پنهان است که در مدل BERT پایه و ParsBERT، ۷۶۸ در نظر گرفته شده است. احتمال آنکه کلمه‌ی i ام که در دنباله‌ی ورودی است، شاخص شروع پاسخ ما باشد، یعنی P_i با فرمول (۵-۱) محاسبه می‌شود:

$$P_i = \frac{e^{S.T_i}}{\sum_j e^{S.T_j}} \quad (۵-۱)$$

در این فرمول، T_i برابر است با بردار لایه‌ی نهایی پنهان برای i امین نشانه که برای BERT پایه و $ParsBERT$ اندازه‌ای برابر با ۷۶۸ دارد. S نیز همان بردار شروع پاسخ است که چند خط قبل از آن صحبت شد. به این ترتیب برای هر نشانه‌ی موجود در متن یک احتمال برای اینکه آن کلمه شروع پاسخ سوال ما باشد، به دست می‌آید.

فرمولی مشابه با رابطه‌ی بالا برای احتمال آنکه توکن i ، توکن پایان پاسخ ما باشد نیز ارائه می‌شود که یک شرط اضافه‌تر دارد. و آن اینست که شاخص پایان جواب باید از شاخص شروع جواب بیشتر باشد. به این معنی که پایان جواب نمی‌تواند پیش از شروع جواب در متن ظاهر شده باشد. رابطه‌ی (۲-۵) این فرمول را نشان می‌دهد.

$$P_i = \frac{e^{E \cdot T_i}}{\sum_j e^{E \cdot T_j}} \quad (۲-۵)$$

در این فرمول، E همان بردار پایان جواب است که در فرآیند تنظیم دقیق معرفی و آموخته می‌شود.

نهایتاً اگر i را شاخص شروع پاسخ و j را شاخص پایان پاسخ در نظر بگیریم، i و j ی که $S \cdot T_i + E \cdot T_j$ را بیشینه می‌کنند به شرطی که $j \geq i$ ، حدس نهایی ما خواهد بود.

با در نظر گرفتن سوالات بدون جواب، فرض می‌کنیم که شاخص شروع و پایان جواب این نوع سوالات، مربوط به نشانه‌ی [CLS] باشد؛ بدین ترتیب همانطور که در بخش قبل برای تک‌تک نشانه‌ها احتمالی برای شروع و پایان پاسخ محاسبه کردیم، این بار باید توکن [CLS] را نیز در محاسباتمان دخیل کنیم. امتیاز مربوط به حالت بدون پاسخ (s_{null}) با فرمول (۳-۵) محاسبه می‌شود که در آن C مربوط می‌شود به بردار نهایی متناظر با [CLS].

$$s_{null} = S \cdot C + E \cdot C \quad (۳-۵)$$

برای محاسبه‌ی امتیاز برای سوالات دارای جواب، رابطه‌ی (۴-۵) را به دست آوردیم.

$$s_{i,j} = \max_{j \geq i} S \cdot T_i + E \cdot T_j \quad (۴-۵)$$

در این صورت زمانی نتیجه می‌گیریم که پرسش داده شده با توجه به متن قابل پاسخ نیست که رابطه‌ی (۵-۵) برقرار باشد.

$$s_{i,j} > s_{null} + \tau \quad (۵-۵)$$

که در آن τ آستانه‌ایست که برای آن که سوالی را بدون جواب تلقی کنیم در نظر گرفته می‌شود.

۵-۲ پیاده‌سازی

برای پیاده‌سازی این پروژه، از اکوسیستم Hugging Face استفاده شده است. Hugging Face در اصل یک شرکت هوش مصنوعی است که اکوسیستمی فراهم آورده که با تجمع و توسعه‌ی کاربردی‌ترین ابزارهای پردازش زبان طبیعی، حل مسائل مربوط به این حوزه را تسهیل بخشیده است. کتابخانه‌های این اکوسیستم، از جمله transformers، یکی از کتابخانه‌های محبوب زبان پایتون به حساب می‌رود.

۵-۲-۱ آماده‌سازی داده‌ها

برای طراحی سیستم پرسش و پاسخ، در ابتدا نیاز است داده‌ها را آماده کنیم. برای آماده‌سازی داده‌های PersianQA از کتابخانه‌ی datasets، استفاده کرده‌ایم. تابع load_dataset به ما کمک می‌کند مجموعه‌داده‌ی پردازش زبان طبیعی مورد نظر را بارگیری و از آن استفاده کنیم. این مجموعه‌داده می‌تواند در هر جایی ذخیره شده باشد. روی Hugging Face Hub، روی دیسک محلی یا بر روی Github. در این پروژه PersianQA را از Hugging Face Hub بارگیری نموده‌ایم.

با استفاده از تابع load_metric می‌توان ملاک‌های ارزیابی مورد نظر را بارگیری کرد. در صورتی که ملاک‌های ارزیابی توسط کتابخانه‌ی datasets به صورت مستقیم، پشتیبانی نشود نیز می‌توان ملاک‌های ارزیابی جدیدی تعریف و از آن‌ها استفاده کرد. معیارهایی که پیش‌تر معرفی کردیم، برای این پروژه توسط کتابخانه‌ی datasets پشتیبانی می‌شود.

پیش از آنکه داده‌ها را وارد مدل کنیم، نیاز است آن‌ها را به نشانه تبدیل کنیم. این عملیات توسط کلاس AutoTokenizer و با تابع from_pretrained قابل انجام است. با این تابع، می‌توانیم Tokenizer از پیش آموزش‌دیده‌ی مد نظر را بارگیری و از آن استفاده کنیم. مدل ParsBERT که می‌خواهیم از آن استفاده کنیم، با آدرس HooshvareLab/bert-fa-zwnj-base از کتابخانه‌ی transformers قابل دسترسی است. Tokenizer آن نیز بر مبنای الگوریتم WordPiece می‌باشد که پیش‌تر به آن پرداختیم.

۵-۲-۲ تنظیم دقیق مدل

برای ساخت مدلی که لایه‌ی خروجی آن، متناسب با یک کار پایین‌دستی خاص است، می‌توان از Automodelهایی که توسط کتابخانه‌ی transformers معرفی می‌شوند استفاده کرد. هر کار پایین‌دستی که بخواهیم انجام بدهیم، AutoModel مخصوص خود را دارد که لایه‌ی آخر آن، به فراخور کار مورد نظر طراحی شده است. در مورد کار پرسش و پاسخ، از AutoModelForQuestionAnswering استفاده کرده‌ایم. با استفاده از متد from_pretrained از این کلاس، می‌توانیم مدلی که می‌خواهیم از آن یک مدل پرسش و پاسخ بسازیم را انتخاب کنیم و وزن‌های آن را بارگیری نماییم. این کلاس، با جایگزین کردن لایه‌ی خروجی مدل از پیش آموزش‌دیده با لایه‌ی جدید متناسب با کار پرسش و پاسخ، مقدمات تنظیم

دقیق را برای ما فراهم می آورد. در این مرحله به تعداد ایپاکی که می خواهیم، مدل را آموزش می دهیم تا وزن های جدید، به دست آیند.

کار آموزش مدل از طریق کلاس Trainer از کتابخانه ی transformers، قابل انجام است. برای آنکه بتوانیم از این کلاس استفاده کنیم، ابتدا باید آرگومان های مربوط به آموزش را تنظیم کرده و به عنوان پارامتر ورودی به تابع train از این کلاس بدهیم. پس از آن با فراخوانی این تابع، فرآیند آموزش شروع می شود. در طول آموزش، پس از پایان هر ایپاک، می توان تابع هزینه ی مربوط به داده های آموزش و داده های اعتبارسنجی یا آزمایش را دید.

۵-۲-۳ ارزیابی مدل

خروجی مدل برای هر نمونه، مقدار تابع هزینه و بردار نهایی متناظر با احتمال آغاز و پایان جواب هر نشانه را بر می گرداند. برای انتخاب بهترین شاخص شروع و پایان، می توان نشانه ای که بهترین امتیاز را برای شاخص شروع بودن دارد، و نشانه ای که بهترین امتیاز را برای شاخص پایان بودن دارد، برداریم و به عنوان جواب در نظر بگیریم. اما همانطور که اشاره کردیم، باید این شرط را در نظر بگیریم که شاخص پایان کوچک تر از شاخص شروع نباشد؛ به همین دلیل ما به تعداد n بهترین کاندید جواب را انتخاب می کنیم. میزان n متغیر است. در این پروژه این عدد برابر ۲۰ در نظر گرفته شده است. از بهترین جواب شروع به بررسی می کنیم و بهترین جوابی که شرط را ارضا کرد را به عنوان جواب نهایی در نظر می گیریم. برای این کار نیازمند یک حلقه ی تودرتو هستیم.

پس از پایان کار مدل، می توان آن را در اکوسیستم Hugging Face قرار داد. این اکوسیستم هم چنین یک رابط کاربری نیز فراهم کرده است که به راحتی می توان به تست مدل پرداخت. تصویر زیر مربوط به مدلیست که تنظیم کرده ایم. این مدل، برای این سوال که «من به چه کاری علاقمندم؟» با توجه به متن «من به کار پردازش زبان های طبیعی علاقه دارم»، جواب «کار پردازش زبان های طبیعی» را برگردانده است.



شکل (۵-۱) تست نمونه ای از پرسش و متن در مدل تنظیم شده با استفاده از Hugging Face API

۳-۵ نتایج ارزیابی

نتایج مربوط به تنظیم ابرپارامترهای نرخ یادگیری، اندازه‌ی دسته و تعداد اپاک این سیستم در یازده حالت و دو ملاک ارزشیابی ما یعنی تطابق دقیق و امتیاز F1 در جدول (۱-۵) آمده است.

امتیاز F1	امتیاز تطابق دقیق	اندازه‌ی دسته	نرخ یادگیری	تعداد اپاک	شماره حالت
53.36	39.56	1	1e-5	2	#1
55.47	40.64	1	1e-5	3	#2
56.12	41.29	1	1e-5	4	#3
56.65	42.79	1	2e-5	2	#4
58.66	43.11	1	2e-5	3	#5
57.13	41.61	1	2e-5	4	#6
57.86	42.90	2	2e-5	2	#7
58.70	42.90	2	2e-5	3	#8
59.12	42.79	2	2e-5	4	#9
58.03	42.15	2	2e-5	5	#10
55.30	40.64	3	2e-5	3	#11

جدول (۱-۵) نتایج حاصل از تنظیم دقیق ParsBERT با مجموعه داده‌ی PersianQA

نتایج تابع هزینه‌ی آموزش و اعتبارسنجی برای هر حالت، در ادامه آمده است.

نرخ یادگیری = 1e-5، اندازه‌ی دسته = 1، تعداد اپاک = 2		
Epoch	Training Loss	Validation Loss
1	3.035800	2.606466
2	2.274800	2.577513

جدول (۲-۵) نتایج حالت اول آزمایش

نرخ یادگیری = 1e-5، اندازه‌ی دسته = 1، تعداد اپاک = 3		
Epoch	Training Loss	Validation Loss
1	3.056700	2.791818
2	2.251500	2.666221
3	1.862000	3.047574

جدول (۳-۵) نتایج حالت دوم آزمایش

نرخ یادگیری = 1e-5، اندازه‌ی دسته = 1، تعداد اپاک = 4		
Epoch	Training Loss	Validation Loss
1	3.051900	2.693343
2	2.259600	2.569704
3	1.841600	2.997958
4	1.182800	3.624714

جدول (۴-۵) نتایج حالت سوم آزمایش

نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 1، تعداد اپاک = 2		
Epoch	Training Loss	Validation Loss
1	3.069700	2.635133
2	2.134000	2.600149

جدول (۵-۵) نتایج حالت چهارم آزمایش

نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 1، تعداد اپاک = 3		
Epoch	Training Loss	Validation Loss
1	3.050500	2.811161
2	2.215500	2.683151
3	1.470300	3.428255

جدول (۶-۵) نتایج حالت پنجم آزمایش

نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 1، تعداد اپاک = 4		
Epoch	Training Loss	Validation Loss
1	3.086400	2.852875
2	2.248100	2.842860
3	1.566200	3.428941
4	0.855400	4.465066

جدول (۷-۵) نتایج حالت ششم آزمایش

نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 2، تعداد اپاک = 2		
Epoch	Training Loss	Validation Loss
1	2.806900	2.307975
2	1.892500	2.277155

جدول (۸-۵) نتایج حالت هفتم آزمایش

نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 2، تعداد اپاک = 3		
Epoch	Training Loss	Validation Loss
1	2.785900	2.301585
2	1.907500	2.281873
3	1.200700	2.742789

جدول (۹-۵) نتایج حالت هشتم آزمایش

نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 2، تعداد اپاک = 4		
Epoch	Training Loss	Validation Loss
1	2.769700	2.342861
2	1.930000	2.334564
3	1.241200	2.754138
4	0.777500	3.464884

جدول (۵-۱۰) نتایج حالت نهم آزمایش

نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 2، تعداد اپاک = 5		
Epoch	Training Loss	Validation Loss
1	2.789600	2.297174
2	1.985400	2.319313
3	1.288500	2.708374
4	0.777500	3.464884
5	0.802800	3.723090

جدول (۵-۱۱) نتایج حالت دهم آزمایش

نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 3، تعداد اپاک = 3		
Epoch	Training Loss	Validation Loss
1	2.758700	2.236697
2	1.926900	2.259194
3	1.247600	2.479626

جدول (۵-۱۲) نتایج حالت یازدهم آزمایش

همانطور که مشاهده می‌شود، بهترین نتیجه با ابرپارامترهای: (نرخ یادگیری = $2e-5$ ، اندازه‌ی دسته = 2، تعداد اپاک = 4) در حالت نهم، حاصل شد.

این آزمایش‌ها بر روی سیستمی با کارت گرافیک NVIDIA GEFORCE GTX 1060M انجام گرفتند. بر روی این سخت‌افزار، هر اپاک در حدود سی دقیقه زمان برد. باید گفت افزایش اندازه‌ی دسته به بیش از سه، ممکن نبود و سیستم را با خطای حافظه‌ی GPU مواجه می‌کرد.

نتیجه گیری

در این پروژه سعی کردیم به مسئله‌ی درک مطلب بپردازیم. ابتدا با لزوم توجه به این قضیه آشنا شدیم. پس از معرفی مجموعه داده‌های مورد استفاده، گفتیم که چرا بهتر است به سراغ معماری‌های غیربازگشتی برویم. در ادامه به مبدل‌ها پرداختیم و بعد یکی از مهم‌ترین مدل‌های مبتنی بر مبدل را مطالعه کردیم. دیدیم که چطور یادگیری انتقالی در بحث پردازش زبان طبیعی به کار می‌رود. در آخر مدل مبتنی بر مبدل ParsBERT را برای کار مورد نظر خود در این پروژه، یعنی پرسش و پاسخ فارسی، تنظیم دقیق کردیم و نتایج حاصل از آن را بررسی نمودیم.

در ادامه‌ی کار این پروژه، می‌توان به جای مدل BERT به سراغ مدل XLNet رفت. XLNet مدلی شبیه به BERT است. با اینهمه این مدل در بیست کار مربوط به پردازش زبان طبیعی از جمله پرسش و پاسخ، بهتر از BERT عمل می‌کند. اگر چه به جهت سنگینی محاسبات، پیچیده‌تر از BERT بوده و زمان بیشتری برای آموزش نیاز دارد. XLNet از قابلیت پردازش متون با طول متغیر برخوردار است؛ هم‌چنین از روشی استفاده می‌کند که در آن، در مرحله‌ی پیش‌آموزش، نیاز به پنهان کردن کلمات ورودی نباشد. این موضوع عملکرد بهتر این مدل را سبب شده است. امید است با پیاده‌سازی سیستم پرسش و پاسخ مبتنی بر XLNet، نتایج ارزیابی، بهبود یابند.

- [1] Amidi, A., And Amidi, S. *Recurrent nueral networks cheatsheet* Available at: <https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks>
- [2] Ayoubi, S., And Davoodeh, M. *PersianQA: a dataset for Persian Question Answering* Available at: <https://github.com/sajjjadayobi/PersianQA>
- [3] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- [4] Farahani, M., Gharachorloo, M., Farahani, M., & Manthouri, M. (2021). Parsbert: Transformer-based model for persian language understanding. *Neural Processing Letters*, 53(6), 3831-3847.
- [5] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735-1780.
- [6] Jay Alammar, (2018, June 27). The Illustrated Transformer [Blog post]. Retrieved from <https://jalammar.github.io/illustrated-transformer/>
- [7] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- [8] Olah, C. Understanding lstm networks. Available at <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- [9] Pennington, J., Socher, R., & Manning, C. D. (2014, October). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (pp. 1532-1543).
- [10] Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (1802). Deep contextualized word representations. *arXiv 2018. arXiv preprint arXiv:1802.05365*, 12.
- [11] Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding with unsupervised learning.

- [12] Rajpurkar, P., Jia, R., & Liang, P. (2018). Know what you don't know: Unanswerable questions for SQuAD. *arXiv preprint arXiv:1806.03822*.
- [13] Rajpurkar, P., Zhang, J., Lopyrev, K., & Liang, P. (2016). SQuAD: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- [14] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.