

# مجموعه داده‌ی چند سطحی فارسی برای بازیابی اطلاعات

علی عابدزاده

کارشناسی ارشد

دانشکده مهندسی کامپیوتر، دانشگاه اصفهان، اصفهان، ایران

رضا رضانی\*

دکتری تخصصی

استادیار، دانشکده مهندسی کامپیوتر، دانشگاه اصفهان، اصفهان، ایران

افسانه فاطمی

دکتری تخصصی

دانشیار، دانشکده مهندسی کامپیوتر، دانشگاه اصفهان، اصفهان، ایران

دریافت: ۱۴۰۱/۱۲/۲۷ | پذیرش: ۱۴۰۲/۰۸/۲۷ | مقاله برای اصلاح به مدت ۵ ماه و ۱۶ روز نزد پدیدآوران بوده است.

نشریه علمی (رتبه بین المللی)  
پژوهشگاه علوم و فناوری اطلاعات ایران  
شاپا(چاپی) ۲۲۵۱-۸۲۲۳  
شاپا(الکترونیکی) ۲۲۵۱-۸۲۳۱  
نمایه در SCOPUS، LISTA و ISC  
http://jipm.irandoc.ac.ir  
دوره XX | شماره X | صص XX-XX  
۱۳XX X

نوع مقاله: پژوهشی

به این مقاله به شکل زیر استناد کنید:

درون متن:

(عابدزاده، رضانی، فاطمی، زودآی‌ند)

در فهرست منابع:

عابدزاده، علی، رضانی، رضا، فاطمی، افسانه،  
زودآی‌ند، مجموعه داده‌ی چند سطحی فارسی برای  
بازیابی اطلاعات، پژوهشنامه پردازش و مدی‌ری‌ت  
(اطلاعات)

(دسترسی در

<http://jipm.irandoc.ac.ir>

**چکیده:** یک سامانه‌ی بازیابی اطلاعات وظیفه دارد با دریافت یک پرسش یا پرسه<sup>۱</sup>، اسناد مرتبط با آن پرسه<sup>۲</sup> را بازیابی کند؛ که این بازیابی از میان مجموعه‌ای بزرگ از هزاران تا میلیون‌ها سند انجام می‌شود. در سال‌های اخیر، پژوهش‌های زیادی برای توسعه‌ی سامانه‌های بازیابی اطلاعات با استفاده از مدل‌های زبان<sup>۳</sup> انجام شده است؛ اما در این زمینه، پژوهشی برای زبان فارسی یافت نشد. یکی از علت‌های اصلی این امر، وجود نداشتن یک مجموعه داده‌ی فارسی مناسب برای آموزش مدل‌های زبان می‌باشد. در این پژوهش، ابتدا یک مجموعه داده‌ی بازیابی اطلاعات فارسی ارائه شده است. پس از آن، روش‌هایی برای غنی‌سازی این مجموعه‌ی داده مورد بحث قرار گرفته است. این غنی‌سازی با کمک چند سطحی کردن ارتباط میان پرسه و سند انجام می‌شود؛ به نحوی که مجموعه داده‌ی جدید می‌تواند رابطه بین پرسه و سند را بجای دو سطح (کاملاً نامرتب - کاملاً مرتب) در چهار سطح (نامرتب - مرتب - بسیار مرتب - کاملاً مرتب) نشان دهد. مجموعه داده ایجاد شده PersianMLIR<sup>۴</sup> نام دارد<sup>۵</sup>. آزمایش‌ها بیانگر بهبود عملکرد سامانه هم برای زبان فارسی و هم برای زبان انگلیسی است که این میزان بهبود برای زبان فارسی ۱۷٪ می‌باشد.

<sup>۲</sup> در این مقاله، به منظور یکنواختی، از کلمه پرسه بجای پرسش یا پرسه استفاده می‌شود

<sup>۳</sup> Language Models

<sup>۴</sup> Persian Multi-Level Information Retrieval (PersianMLIR)

<sup>۵</sup> لینک دانلود مجموعه داده: <https://github.com/BigData-IsfahanUni/PersianMLIR>

کلیدواژه ها: بازیابی اطلاعات، مدل های زبان، مجموعه داده بازیابی اطلاعات،  
مجموعه داده فارسی

\* رضا رضانی r.ramezani@eng.ui.ac.ir

#### ۱. مقدمه

سامانه های بازیابی اطلاعات متنی کاربردهای بسیاری دارند. طی دهه ها پژوهش در این زمینه، محققان روش های متفاوتی برای ساخت سامانه های بازیابی اطلاعات ارائه کرده اند. از مشهورترین و پر استفاده ترین روش های بازیابی اطلاعات می توان به TF-IDF (Salton and Buckley 1988) و BM25 (Robertson and Zaragoza 2009) اشاره کرد. این دو روش، بر اساس مدل بسته کلمات و شباهت لغوی عمل می کنند. بر اساس معیارهای Trec-DL (Craswell et al. 2020; Craswell, ) و (Mitra, Yilmaz, and Campos 2021)، سامانه های بازیابی اطلاعات را بر اساس راه کار شباهت سنجی، می توان به سه دسته تقسیم کرد:

- سامانه های سنتی: این سامانه ها بر اساس روش های قدیمی مانند BM25 و TF-IDF (و به طور کلی، شباهت لغوی میان اسناد و پرسه) کار می کنند. در این سامانه ها استفاده از نمایش  $\chi^2$  اسناد بسیار رایج می باشد.
- سامانه های شبکه عصبی: در این سامانه ها، برای بازیابی از شبکه های عصبی کمک گرفته می شود. تمام سامانه هایی که بر اساس شبکه های عصبی هستند، اما شبکه های عصبی مورد استفاده، یک مدل زبان نمی باشد در این دسته قرار می گیرند. در این سامانه ها استفاده از نمایش متراکم اسناد بسیار رایج می باشد.
- سامانه های شبکه عصبی مدل زبان: اگر شبکه عصبی مورد استفاده در یک سامانه از نوع مدل زبان (مانند BERT (Devlin et al. 2019)) باشد، آن سامانه در این دسته قرار می گیرد. در این سامانه ها نیز، استفاده از نمایش متراکم اسناد بسیار رایج می باشد.

از اصلی‌ترین اشکالات روش‌های سنتی مبتنی بر نمایش تُنک اسناد، بوجود آمدن اختلاف واژگانی میان سند و پرسه کاربر می‌باشد (برای نمونه «سرخ» و «قرمز») (Mitra and Craswell 2018; Qu et al. 2021). اشکال دیگر نمایش تُنک اسناد زمانی پدیدار می‌شود که یک کلمه دو یا چند مفهوم متفاوت داشته باشد (مانند «شیر»). در نحوه‌ی نمایش متراکم اسناد که معمولاً با کمک مدل‌های زبان و شبکه‌های عصبی تولید می‌شوند، معایب روش‌های قبل تا حد زیادی حل شده است. در این روش‌ها، به جای توجه به تشابه لغوی، تلاش می‌شود تا معنا و مفهوم پرسه و سند درک شود. در این حالت، نه تنها مشکل اختلاف واژگانی حل می‌شود، بلکه یک بازیابی معنایی وجود خواهد داشت (Z. Liu et al. 2021; Qu et al. 2021). اگرچه این روش‌ها دقت بالاتری دارند، اما در مواجهه با واژگان نادر و کمیاب، عملکرد خوبی ندارند و در چنین مواردی، نحوه‌ی نمایش تُنک اسناد می‌تواند نتایج بهتری را حاصل کند (Mitra and Craswell 2018).

در سال‌های اخیر، پژوهش‌های بسیاری برای ساخت سامانه‌های بازیابی اطلاعات با کمک مدل‌های زبان انجام شده است. از آنجایی که مدل‌های زبان عمدتاً بر اساس شبکه‌های عصبی عمیق هستند، وجود مجموعه داده‌های بزرگ مقیاس برای آموزش آن‌ها ضروری می‌باشد (Zhang, Yates, and Lin 2020). در طی سال‌های اخیر، مدل‌های زبان زیادی ساخته شده‌اند و تعداد قابل توجهی از این مدل‌ها قادر به درک و فهم زبان فارسی می‌باشند. با این حال، بدلیل عدم وجود یک مجموعه داده‌ی مناسب، پژوهشی روی بازیابی اطلاعات فارسی با کمک مدل‌های زبان انجام نشده است.

در این پژوهش، ابتدا روشی برای ساخت یک مجموعه داده‌ی بازیابی اطلاعات با کمک مجموعه داده‌های درک مطلب ماشینی<sup>۶</sup> فارسی موجود ارائه شده است.<sup>۷</sup> سپس روشی برای غنی‌سازی این مجموعه داده‌ی جدید ارائه گردیده است. در مجموعه داده تولید شده اولیه، همانند مجموعه داده‌های بازیابی اطلاعات مرسوم، به ازای هر پرسه و سند، تنها یکی از دو رابطه «کاملاً نامرتب» یا «کاملاً مرتب» وجود دارد. هرچند این مجموعه داده تولید شده، اولین مجموعه داده بازیابی اطلاعات برای زبان فارسی می‌باشد، اما در ادامه به منظور بهبود عملکرد آن، سطح روابط میان هر پرسه و سند را از دو سطح به چهار

<sup>۶</sup> Machine Reading Comprehension

<sup>۷</sup> این مجموعه داده درک زبان ماشینی فارسی در خوشه تحقیقاتی کلان داده دانشگاه اصفهان توسعه داده شده است.

سطح افزایش می‌دهیم به نحوی که ارتباط بین پرسه و سند می‌تواند یکی از چهار حالت «نامرتب» و «مرتبط»، «بسیار مرتبط» و «کاملاً مرتبط» باشد. سپس یک سامانه بازیابی اطلاعات مبتنی بر مدل زبان روی هر دو مجموعه داده آموزش داده شده است؛ که نتایج ارزیابی نشان می‌دهد مجموعه داده چهار سطحی می‌تواند بطور میانگین عملکرد مجموعه داده دو سطحی را ۱.۸۷٪ بهبود دهد.

## ۲. پیشینه

از آنجایی که نوآوری اصلی این پژوهش ارائه یک مجموعه داده بازیابی اطلاعات فارسی چند سطحی است، در این بخش به بررسی مجموعه داده‌های بازیابی اطلاعات و درک مطلب ماشینی انگلیسی و فارسی می‌پردازیم. علت بررسی مجموعه داده‌های درک مطلب ماشینی این است که این مجموعه داده‌ها شباهت زیادی به مجموعه داده‌های بازیابی اطلاعات دارند و در مواردی می‌توان با اعمال تغییراتی در آن‌ها، یک مجموعه داده‌ی بازیابی اطلاعات ایجاد کرد. البته تبدیل کردن یک مجموعه داده‌ی درک مطلب ماشینی به یک مجموعه داده‌ی بازیابی اطلاعات می‌تواند مشکلاتی چون سوگیری داده‌های آموزشی، بی معنا بودن پرسه‌ها در صورت وجود نداشتن پاراگراف مربوطه، همپوشانی لغوی زیاد بین پرسه و سند مرتبط را بوجود آورد. با این وجود، تاکنون مجموعه داده‌های بازیابی اطلاعات مختلفی از روی این مجموعه داده‌های درک مطلب ماشینی ایجاد شده است.

مجموعه داده MS Marco (Bajaj et al. 2016; Craswell, Mitra, Yilmaz, Campos, et al. 2021) را می‌توان بزرگ‌ترین مجموعه داده بازیابی اطلاعات انگلیسی دانست که حاوی بیش از ۵۰۰ هزار پرسه و حدود ۸.۸ میلیون سند (پاراگراف) است. این مجموعه داده، یک مجموعه داده دوسطحی محسوب می‌شود و مدل‌های بازیابی اطلاعات بسیاری در مرحله‌ی آموزش از آن استفاده می‌کنند. Trec-DL یکی دیگر از مجموعه داده‌های معروف بازیابی اطلاعات است که در آن برای هر پرسه، تعداد زیادی سند نشانه‌گذاری شده وجود دارد. بر خلاف مجموعه داده‌ی MS Marco که یک مجموعه‌ی دوسطحی می‌باشد، مجموعه داده Trec-DL چندسطحی می‌باشد که سطوح روابط بین پرسه‌ها و اسناد به صورت زیر می‌باشد:

- کاملاً مرتبط (سطح ۳): سند متعلق به پرسه می‌باشد و پاسخ دقیق پرسه در آن قرار دارد.
- بسیار مرتبط (سطح ۲): سند تا حدودی پاسخ پرسه را در خود دارد؛ اما پاسخ دقیق واضح نیست و یا در میان اطلاعات اضافه، پنهان شده است.
- مرتبط (سطح ۱): سند به نظر در ارتباط با پرسه می‌رسد، اما پاسخ در آن وجود ندارد.

- نامرتب (سطح ۰): سند هیچ ارتباطی با پرسه ندارد.

مجموعه داده SQuAD (Rajpurkar et al. 2016; Rajpurkar, Jia, and Liang 2018) یک مجموعه داده درک مطلب ماشینی انگلیسی است که شامل بیش از ۱۰۰ هزار پرسه و پاراگراف می‌باشد. مجموعه داده‌های درک مطلب ماشینی معمولاً یک مجموعه از سه تایی‌ها به فرم (پاراگراف، پرسه، پاسخ پرسه) می‌باشند. با استفاده از این سه تایی‌ها، می‌توان درک مطلب ماشینی را آموزش داد که با دریافت پرسه و پاراگراف، سعی می‌کند تا پاسخ پرسه را از داخل پاراگراف پیدا کند. بطور مشابه، مجموعه داده‌ی بازیابی اطلاعات ANTIQUE (Hashemi et al. 2020) شامل بیش از ۲۵۰۰ پرسه پیچیده و دشوار است که از تالارهای گفت و گوی آنلاین جمع‌آوری شده‌اند. در این مجموعه داده، به ازای هر پرسه، حدود ۱۳ پاراگراف با سطوح ارتباطی مختلف توسط انسان نشانه‌گذاری شده است. مجموعه داده‌ی NaturalQuestions (Kwiatkowski et al. 2019) نیز یک مجموعه داده‌ی نسبتاً بزرگ مقیاس با بیش از ۳۰۰ هزار داده‌ی آموزشی می‌باشد. در این مجموعه داده که برای آموزش و ارزیابی درک مطلب ماشینی طراحی شده است، هر پرسه معمولاً با یک پاسخ کوتاه و یک پاسخ بلند و نیز پاراگراف مربوطه همراه شده است.

مجموعه داده‌هایی که تا کنون مورد بحث قرار گرفتند، مجموعه داده‌های انگلیسی بودند. با بررسی‌های صورت گرفته، مجموعه داده بازیابی اطلاعات فارسی که برای آموزش یک مدل زبان مناسب باشد، شناسایی نشد. مجموعه داده مورد نظر این پژوهش، باید ساختاری شبیه به مجموعه داده‌های انگلیسی بازیابی اطلاعات (که در بالا بررسی شدند) داشته باشد. خوشبختانه در سال‌های اخیر، فعالیت‌های زیادی در زمینه‌ی پرسش و پاسخ فارسی و درک مطلب ماشینی فارسی صورت گرفته و مجموعه داده‌هایی برای این زمینه‌های پژوهشی فراهم شده که می‌توان از آن‌ها به منظور ایجاد مجموعه داده بازیابی اطلاعات فارسی استفاده نمود. به عنوان مثال مجموعه ParsiNLU (Khashabi et al. 2021) متشکل از تعدادی مجموعه داده است که می‌توان از آن برای وظایف مختلف پردازش زبان طبیعی مانند استنتاج متنی، بازنویسی پرسه، پرسش پاسخ چند جوابی، ترجمه ماشینی، تحلیل احساسات و درک مطلب ماشینی استفاده کرد. بخش درک مطلب ماشینی این مجموعه شامل ۱۳۰۰ پرسه می‌باشد که تمام این پرسه‌ها با کمک موتور جست‌وجوی گوگل جمع‌آوری شده است.

مجموعه داده PersianQA (Ayoubi Sajjad & Davoodeh 2021) یک مجموعه داده‌ی فارسی برای درک مطلب ماشینی است که نحوه‌ی ایجاد و ساختار کلی آن مشابه SQuAD می‌باشد. این مجموعه داده شامل حدود ۱۰ هزار داده‌ی آموزشی می‌باشد و اسناد این مجموعه از ویکی‌پدیای فارسی استخراج شده است. بطور مشابه، مجموعه داده ParSQuAD (Abadani et al. 2021) که یک مجموعه‌ی فارسی برای درک مطلب ماشینی است، نسخه‌ی ترجمه شده‌ی SQuAD می‌باشد. مجموعه داده PersianQuAD (Kazemi, Mozafari, and Nematbakhsh 2022) نیز مانند سایر مجموعه‌های فارسی، یک مجموعه داده‌ی درک مطلب ماشینی می‌باشد که فرآیند ساخت این مجموعه داده کاملاً مشابه مجموعه داده‌ی SQuAD بوده و شامل حدود ۲۰ هزار داده‌ی آموزشی شامل سه تایی‌های (پرسه، پاسخ، سند مرتبط) می‌باشد. در این مجموعه نیز از اسناد ویکی‌پدیای فارسی برای طراحی مجموعه استفاده شده و برای طراحی پرسه‌ها و نشانه‌گذاری آن‌ها از افراد فارسی زبان کمک گرفته شده است.

از میان مجموعه داده‌های فارسی نام برده شده، مجموعه‌ی PersianQuAD مناسب‌ترین انتخاب برای ساخت یک مجموعه داده‌ی جدید بازیابی اطلاعات فارسی می‌باشد. علت این امر آن است که برخلاف ParsQuAD (که یک نسخه‌ی ترجمه شده می‌باشد)، اسناد مورد استفاده برای ساخت این مجموعه در دسترس می‌باشند (ویکی‌پدیای فارسی) و این امر برای ساخت یک مجموعه داده‌ی بازیابی اطلاعات ضروری می‌باشد. همچنین این مجموعه داده، تعداد داده‌ی آموزشی بیشتری نسبت به PersianQA و ParsiNLU دارد و تعداد نمونه‌های آموزشی از اهمیت زیادی برخوردار می‌باشد (Zhang et al. 2020). در مورد ساخت مجموعه داده بازیابی اطلاعات با استفاده از این مجموعه داده، در بخش بعد صحبت شده است.

### ۳. ساخت مجموعه داده بازیابی اطلاعات

همانطور که پیش‌تر بیان شد، یکی از نوآوری‌های این پژوهش ارائه یک مجموعه داده بازیابی اطلاعات فارسی می‌باشد. در این بخش ابتدا نحوه‌ی تبدیل مجموعه داده PersianQuAD (که یک مجموعه داده‌ی درک مطلب ماشینی است) به یک مجموعه داده‌ی بازیابی اطلاعات بیان شده است. مجموعه داده‌ی ساخته شده اولیه در این بخش یک مجموعه‌ی دوسطحی است. سپس نحوه غنی‌سازی این مجموعه داده به کمک افزایش سطح روابط از دو سطح به چهار سطح مورد بحث قرار گرفته است.

### ۳-۱ ساخت مجموعه داده‌ی بازیابی اطلاعات فارسی دوسطحی

همان‌طور که اشاره شد، برای زبان فارسی، مجموعه داده‌ی بازیابی اطلاعاتی که برای آموزش و تنظیم دقیق مدل‌های زبان مناسب باشد وجود ندارد، اما با کمک مجموعه داده‌های فارسی درک مطلب ماشینی، می‌توان یک مجموعه داده‌ی جدید تولید کرد. برای این منظور در این پژوهش از مجموعه داده‌ی PersianQuAD (Kazemi et al. 2022) استفاده شده است. در این مجموعه داده، به ازای هر پرسه، چندین سند مرتبط وجود دارد. در این نوشتار، منظور از «سند» یک سند کامل، یک پاراگراف یا حتی یک جمله است.

در این پژوهش از DPR به عنوان مدل بازیابی اطلاعات استفاده خواهد شد (Karpukhin et al. 2020). لذا ساختار مجموعه داده بازیابی اطلاعات فارسی بایستی به نحوی شود که قابل استفاده توسط DPR باشد. DPR یک مدل بازیابی اطلاعات می‌باشد که با کمک آن می‌توان اسناد مرتبط با یک پرسه را بازیابی کرد. در این مدل از دو شبکه‌ی عصبی مجزا برای تعبیه کردن پرسه‌ها و اسناد استفاده می‌شود و شباهت میان خروجی‌های این دو شبکه نمایانگر میزان ارتباط میان پرسه و سند خواهد بود. مجموعه داده‌ی مورد استفاده در DPR شامل سه فایل می‌باشد:

- **فایل مجموعه‌ی اسناد:** یک پیکره متنی بزرگ که بازیابی بر روی آن انجام می‌شود. در این پیکره، هر سند می‌تواند یک سند کامل، یک پاراگراف یا یک جمله باشد.
- **فایل مجموعه‌ی آموزشی:** داده‌های آموزشی مدل که هر نمونه شامل یک پرسه، اسناد مرتبط و همین‌طور اسناد نامرتبط با آن پرسه می‌باشد.
- **فایل مجموعه‌ی ارزیابی:** مانند مجموعه‌ی آموزشی است اما با هدف ارزیابی سامانه استفاده می‌شود.

نحوه‌ی ساخت این سه فایل در ادامه بیان شده است.

### ۳-۱-۱ ساخت مجموعه‌ی اسناد

در ساخت PersianQuAD ابتدا با کمک الگوریتم PageRank حدود ۴۴۰۰ سند ویکی‌پدیا انتخاب شده است (Kazemi et al. 2022). سازندگان این مجموعه داده با استفاده از PageRank اطمینان حاصل کرده‌اند که صفحات مهم و پرارجاع ویکی‌پدیا که شامل اطلاعات داغ و محبوبی هستند انتخاب شده‌اند. سپس با کمک ۱۹۰۵ سند (که می‌توان آن‌ها را به ۲۶ هزار پاراگراف تبدیل کرد)،

مجموعه‌ی داده‌ی درک مطلب ماشینی تولید شده است.<sup>۸</sup> در این پژوهش برای ساختن مجموعه‌ی اسناد، از این ۲۶ هزار پاراگراف استفاده شده است. با این وجود، از آنجایی که قصد داریم مجموعه‌ی اسناد مورد استفاده در این پژوهش، کمی بزرگ مقیاس‌تر از این تعداد سند باشد، لازم است تا اسناد دیگری نیز به این مجموعه اضافه شوند. برای این منظور، ابتدا یک نسخه‌ی متنی کامل از ویکی‌پدیای فارسی بارگیری شده است. پس از آن، متون داخل این فایل استخراج شده و به پاراگراف‌هایی با اندازه‌ی حداکثر ۵۰۰ واژه تبدیل شده‌اند. در نهایت، با انتخاب پاراگراف‌ها بصورت تصادفی، اندازه‌ی مجموعه‌ی اسناد به ۲۰۰ هزار پاراگراف رسیده است. لازم به ذکر است که برای جلوگیری از وارد شدن پاراگراف‌های تکراری در مجموعه‌ی اسناد، از وارد شدن پاراگراف‌هایی که عنوان سند آن‌ها، در لیست اسناد مورد استفاده‌ی PersianQuAD وجود داشته است خودداری شده است.

### ۳-۱-۲ ساخت مجموعه‌ی آموزشی

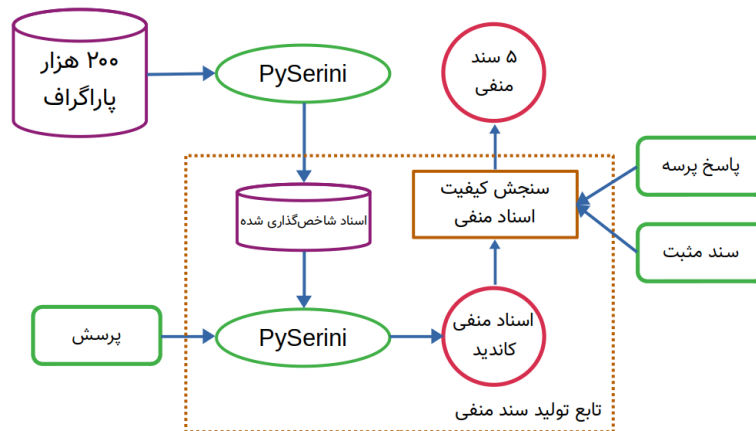
در مجموعه‌داده‌ی آموزشی بازیابی اطلاعات، هر نمونه یک سه‌تایی به فرم (پرسه، مجموعه اسناد مرتبط، مجموعه اسناد نامرتب) می‌باشد. با کمک PersianQuAD می‌توان دو عنصر اول این سه‌تایی را تولید کرد، اما برای تولید مجموعه اسناد نامرتب، باید از مجموعه‌ی کل اسناد (که در مرحله‌ی قبل ساخته شده) کمک گرفته شود. اسناد مرتبط معمولاً نشانه‌گذاری شده‌اند؛ اما برای انتخاب اسناد نامرتب، نیاز به یک مجموعه‌ی اسناد بزرگ است و لذا فرض می‌شود که هر سندی در این مجموعه، به غیر از اسناد مرتبط که نشانه‌گذاری شده‌اند، اسناد نامرتب و منفی هستند (Karpukhin et al. 2020). روش‌های زیادی برای انتخاب سند منفی وجود دارد (Karpukhin et al. 2020; Y. Liu et al. 2021; Qu et al. 2021) اما یکی از رایج‌ترین روش‌ها، استفاده از اسناد با رتبه‌ی BM25 بالا می‌باشد که حاوی پاسخ نمی‌باشند. در این پژوهش برای انتخاب اسناد منفی، ابتدا اسناد با کمک ابزار PySerini (Lin et al. 2021; Yang, Fang, and Lin 2017) شاخص‌گذاری شده‌اند. PySerini ابزاری برای شاخص‌گذاری و بازیابی اطلاعات بر روی اسناد می‌باشد که می‌توان از آن برای بازیابی اسناد با کمک BM25 استفاده کرد. پس از شاخص‌گذاری اسناد، به ازای هر پرسه، پنج سند برتر BM25 که شامل پاسخ نمی‌باشند انتخاب می‌شوند. علاوه بر شرایط ذکر شده، اطمینان حاصل می‌شود که اسناد منفی انتخاب شده، عنوان سند متفاوتی از عنوان سند مرتبط داشته باشند. به این ترتیب، می‌توان اطمینان بیشتری حاصل نمود که اسناد منفی انتخاب شده، ارتباطی با پرسه مد نظر ندارند. در شکل ۱ کلیات فرآیند انتخاب

<sup>۸</sup> این ۱۹۰۵ سند از طریق سازندگان PersianQuAD در اختیار ما قرار گرفته است.



پنج سند منفی به ازای هر پرسه نشان داده شده است. همانطور که در شکل دیده می‌شود، پس از شاخص گذاری اسناد با کمک PySerini، به ازای هر پرسه‌ی ورودی، تعدادی سند بازیابی می‌شود که به عنوان اسناد منفی کاندید در نظر گرفته می‌شوند. پس از آن با کمک محدودیت‌هایی که به آن‌ها پرداخته شد (عدم وجود پاسخ در سند کاندید و تفاوت عنوان سند کاندید با سند مثبت)، این اسناد منفی فیلتر می‌شوند تا اسناد منفی باکیفیت تری تولید شود. اعمال کردن این محدودیت‌ها به عهده‌ی بخش «سنجش کیفیت اسناد منفی» است. از میان اسناد منفی فیلتر شده، پنج سند با بالاترین امتیاز BM25 به عنوان اسناد منفی انتخاب می‌شوند.

با کمک روش بیان شده، می‌توان مجموعه داده‌ی آموزشی مورد نیاز برای آموزش DPR را ایجاد نمود. چگونگی عملکرد و نحوه‌ی آموزش مدل DPR پایه در بخش (۴-۱) توضیح داده شده است.



شکل ۱: فرآیند انتخاب اسناد منفی برای هر پرسه

### ۳-۱-۳ ساخت مجموعه‌ی ارزیابی

روند ساخت مجموعه‌ی ارزیابی بازیابی اطلاعات، بسیار مشابه روند ساخت مجموعه‌ی آموزشی می‌باشد؛ با این تفاوت که در این مجموعه‌ی آموزشی، نیازی به وجود اسناد نامرتب و منفی نمی‌باشد؛ اما در عوض، نیاز به دانستن پاسخ مربوط به پرسه می‌باشد. در این حالت سه تایی‌های (پرسه، پاسخ پرسه، سند مرتبط) تشکیل می‌شوند. استفاده از پاسخ پرسه در سامانه‌های بازیابی اطلاعات امر رایجی نمی‌باشد، اما DPR از این پاسخ‌ها برای محاسبه‌ی سنجه‌ی Exact Match و سنجش عملکرد مدل تحت آموزش

استفاده می‌کند. در بخش (۵-۱) به شرح سنجه‌های مورد استفاده در DPR و این پژوهش پرداخته شده است.

### ۲-۳ ساخت مجموعه داده‌ی بازیابی اطلاعات فارسی چندسطحی

مجموعه داده‌ی ارائه شده در بخش (۳-۱) یک مجموعه داده‌ی دوسطحی است؛ چرا که برای هر پرسه، تنها اسناد «مرتبط» و «نامرتب» تعیین شده است. یک مجموعه داده‌ی چندسطحی می‌تواند رابطه‌ی یک پرسه و سند را در چند سطح نمایش دهند. در این پژوهش از ساختار مجموعه داده Trec-DL که یک مجموعه داده‌ی بازیابی اطلاعات چندسطحی است (Craswell et al. 2020; Craswell, Mitra, Yilmaz, and Campos 2021) استفاده شده است.

مجموعه داده Trec-DL روابط را در چهار سطح نشان می‌دهد: کاملاً مرتبط (سطح ۳)، بسیار مرتبط (سطح ۲)، مرتبط (سطح ۱)، نامرتب (سطح ۰). برای نمونه، اگر پرسه‌ی مورد بررسی «مصدق در چه سالی متولد شده؟» باشد، پاراگراف «مصدق در سال ۱۲۶۱ چشم به جهان گشود و... کاملاً مرتبط، پاراگراف «در سال ۱۲۸۰ خورشیدی، مصدق که در آن زمان نوزده سال داشت... بسیار مرتبط، پاراگراف «... ازدواج این دو ۶۴ سال تا پایان زندگانی‌شان ادامه یافت» مرتبط و پاراگراف «پس از مرگ مظفرالدین شاه پسرش محمد علی شاه در تاریخ ۲۹ دی ۱۲۸۵ تاج‌گذاری کرد.» نامرتب می‌باشد. در واقع می‌توان گفت که سطح ۰ همان سند منفی است و سطح ۳ نیز همان سند مثبت است. در این حالت، اطلاعات اضافه‌ای میان این دو سطح اضافه شده است (سطح ۱ و ۲) که مدل می‌تواند با کمک آن‌ها، بازیابی دقیق‌تر و بهتری انجام دهد.

در این پژوهش، برای تبدیل کردن مجموعه داده‌ی دوسطحی ساخته شده در بخش (۳-۱) به یک مجموعه‌ی چندسطحی، از دو روش استفاده شده است:

- روش مبتنی بر قوانین دست‌ساز
- روش مبتنی بر طبقه‌بندی و ویژگی‌های لغوی

### ۱-۲-۳ روش مبتنی بر قوانین دست‌ساز

برای ساخت یک مجموعه داده‌ی چند سطحی با کمک قوانین دست‌ساز، یکصد سند برتر مرتبط با هر پرسه با کمک PySerini استخراج شده‌اند تا بتوان از آن‌ها به عنوان اسناد کاندید<sup>۹</sup> استفاده کرد. پس

<sup>۹</sup> سند کاندید سندی است که از میان کل اسناد انتخاب شده است اما میزان ارتباط آن با پرسه‌ی مورد بررسی مشخص نمی‌باشد و می‌تواند به صورت بالقوه در یکی از چهار سطح ارتباطی قرار بگیرد.

از آن، با کمک قوانینی که به صورت دستی تعریف شده‌اند، هر سند کاندید به یکی از سطوح ۰ الی ۲ منتصب شده است (نامرتب، مرتب، بسیار مرتب). سطح ۳ (کاملاً مرتب) تنها شامل یک سند می‌باشد و آن سند همان سند مثبت می‌باشد و قوانین مورد بحث نقشی در انتخاب آن ندارند. برای تعیین سطح یک سند کاندید، علاوه بر خود سند کاندید و عنوان آن، سند مثبت و عنوان آن و همینطور پاسخ کوتاه و ریشه‌های آن مورد نیاز هستند. برای استخراج ریشه‌های پاسخ از ابزار Hazm استفاده شده است<sup>۱۱</sup>.  
قوانین تعریف شده به ترتیب زیر می‌باشند:

- قانون ۱: اگر عنوان سند کاندید و عنوان سند مثبت یکی باشد و همچنین پاراگراف کاندید شامل بخشی از ریشه لغت‌های پاسخ باشد، سند کاندید یک سند سطح ۲ خواهد بود.
  - قانون ۲: اگر پاسخ عیناً در سند کاندید وجود داشته باشد، سند کاندید یک سند سطح ۲ خواهد بود.
  - قانون ۳: اگر عنوان سند کاندید و عنوان سند مثبت یکی باشد یا پاراگراف کاندید شامل بخشی از ریشه لغت‌های پاسخ باشد، سند کاندید یک سند سطح ۱ خواهد بود.
  - قانون ۴: اگر عنوان سند کاندید و عنوان سند مثبت متفاوت باشد و همینطور سند کاندید شامل پاسخ دقیق نباشد، سند کاندید یک سند سطح ۰ خواهد بود.
- از میان قوانین تعریف شده، اولین قانونی که برقرار شود، تعیین‌گر سطح سند می‌باشد و قوانین بعد از آن اثری نخواهند داشت.

### ۲-۲-۳ روش مبتنی بر طبقه‌بندی و ویژگی‌های لغوی

ایده‌ی اصلی این بخش این است که با کمک یک مجموعه داده‌ی چندسطحی انگلیسی موجود، یک مدل آموزش داده شود که بتواند با کمک این مجموعه‌ی داده، میزان ارتباط اسناد را یاد بگیرد. در واقع هدف آن است تا دانش موجود در یک مجموعه داده‌ی چندسطحی به یک مدل منتقل شود و با کمک آن مدل، یک مجموعه داده‌ی جدید ایجاد شود. برای این منظور در این پژوهش، از مجموعه داده‌ی Trec-DL (که یک مجموعه داده‌ی چندسطحی انگلیسی است) برای آموزش مدل استفاده گردیده و از

<sup>۱۱</sup> ابزاری است مانند NLTK پایتون برای تمیز کردن متن، تجزیه نحوی، ریشه‌یابی و ... - [https://www.roshan-](https://www.roshan-ai.ir/hazm)

مدل ایجاد شده، برای ساخت یک مجموعه داده‌ی فارسی استفاده شده است. برای آن منظور، ابتدا نیاز به یک تابع برای استخراج ویژگی می‌باشد، که این ویژگی‌ها عبارتند از:

- آیا سند کاندید شامل پاسخ کامل می‌باشد؟ به عبارتی آیا پاسخ ارائه شده برای پرسه، عیناً در سند کاندید وجود دارد.

- آیا سند کاندید و سند مثبت، عنوان مشترک دارند؟

- لیستی از شباهت‌های لغوی میان سند کاندید و سند مثبت. در این لیست، با کمک سه تابع شباهت کسینوسی، جاکارد و اقلیدوسی، شش مقدار مشابهت میان سند کاندید و سند مثبت محاسبه می‌شوند (دو مقدار به ازای هر تابع شباهت)<sup>۱۱</sup>.

- لیستی از شباهت معنایی میان سند کاندید و سند مثبت. در این لیست، از مدل زبان‌های تولید شده توسط DPR دوسطحی استفاده می‌شود. برای این منظور، DPR به تعداد یک، سه و شش دور تنظیم دقیق می‌شود و از مدل‌های هر اجرا، برای تعبیه کردن و به دست آوردن شباهت میان بردار سند مثبت و بردار سند کاندید استفاده می‌شود. به طور مشابه، میزان شباهت پرسه و سند کاندید نیز با کمک این مدل‌ها محاسبه می‌شود. علاوه بر سه مدل تنظیم دقیق شده‌ی بالا، یک مدل BERT نیز (بدون تنظیم دقیق شدن) برای تعبیه کردن اسناد و پرسه‌ها استفاده می‌شود. برای به دست آوردن شباهت، از توابع ضرب داخلی و همینطور شباهت کسینوسی استفاده شده است. حاصل تمام این محاسبات، ۱۶ مقدار مشابهت (دو تابع شباهت سنجی اعمال شده بر روی هشت مدل زبان) می‌شود که در یک لیست قرار می‌گیرند.

با در دست داشتن این تابع، می‌توان با کمک PySerini یکصد سند کاندید را استخراج و ویژگی‌های هر یک از آنان را استخراج کرد. سپس با وارد کردن این ویژگی‌ها به یک طبقه‌بند می‌توان سطح این اسناد کاندید را تعیین کرد. مانند قبل، تعیین سطح سند تنها برای سطوح ۰ الی ۲ می‌باشد و مجموعه اسناد سطح ۳ مرتبط با هر پرسه تنها شامل یک عضو است که همان سند مثبت می‌باشد و طبقه‌بند نقشی در تعیین آن ندارد.

<sup>۱۱</sup> برای محاسبه‌ی این شباهت‌ها، ابتدا یک بردار از سند کاندید و یک بردار از سند مثبت ساخته می‌شود. هر عنصر در این دو بردار نشان‌دهنده‌ی یک واژه می‌باشد. مقداردهی این عناصر می‌تواند به صورت دودویی (۱ برای وقوع واژه در سند و ۰ برای عدم وقوع آن) و غیر دودویی (تعداد وقوع واژه) باشد. در نتیجه به ازای هر سند یک بردار دودویی و یک بردار غیر دودویی وجود خواهد داشت و با کمک سه تابع شباهت یاد شده، می‌توان شش مقدار شباهت تولید کرد.

مدلهایی که در این بخش برای طبقه‌بندی اسناد مورد استفاده قرار گرفته‌اند، عبارتند از: درخت تصمیم، جنگل تصادفی، ماشین بردار پشتیبان، شبکه‌ی عصبی، نزدیک‌ترین همسایه،<sup>۱۲</sup> AdaBoost،<sup>۱۳</sup> SGD

برای آموزش این مدل‌ها، از ویژگی‌های استخراج شده در مراحل قبل استفاده شده است و برای سطح ارتباط که بایستی پیش‌بینی شود نیز از سطوح تعیین شده در Trec-DL (که توسط سازندگان مجموعه‌ی داده برای هر پاراگراف مشخص شده است) استفاده گردیده است. در میان این ویژگی‌ها، مواردی که از پاسخ استفاده می‌کنند وجود ندارند؛ چرا که مجموعه‌داده‌ی Trec-DL چنین اطلاعاتی را در خود ندارد. از آنجایی که در میان این ویژگی‌ها، ویژگی‌های مرتبط با پاسخ وجود ندارد، سایر ویژگی‌های باقی‌مانده مستقل از زبان می‌باشند و می‌توان توابع مورد استفاده برای مجموعه‌ی انگلیسی را برای استخراج ویژگی از مجموعه‌ی فارسی استفاده کرد. پس از آموزش و ارزیابی مدل، می‌توان مانند بخش قبل، اسناد کاندید فارسی را با کمک PySerini و BM25 تولید نموده، ویژگی‌های آن‌ها را استخراج کرد و سپس آن‌ها را طبقه‌بندی نمود. در نهایت با کمک طبقه‌بندی‌های انجام شده، می‌توان یک مجموعه‌داده‌ی چندسطحی فارسی ساخت.

#### ۴. آموزش مدل بازیابی اطلاعات

##### ۴-۱ مدل DPR پایه

برای استفاده از مجموعه‌داده بازیابی اطلاعات دوسطحی ساخته شده، یک مدل بازیابی اطلاعات به عنوان مدل پایه مورد نیاز است. بدین منظور در این پژوهش از مدل DPR<sup>۱۴</sup> (Karpukhin et al. 2020) استفاده شده است. DPR یک مدل بازیابی اطلاعات می‌باشد که با کمک مدل‌های زبان قادر است تا اسناد مرتبط با یک پرسه را بازیابی کند. DPR از دو رمزگذار BERT مستقل استفاده می‌کند، اما نیازی به پیش‌آموزش آن‌ها از قبل ندارد. در عوض، DPR بر یادگیری یک مدل بازیابی اطلاعات قوی با استفاده از پرسش‌ها و پاسخ‌های دوتایی تمرکز می‌کند. در واقع DPR راه‌هایی را برای انتخاب و

<sup>۱۲</sup> Adaptive Boosting

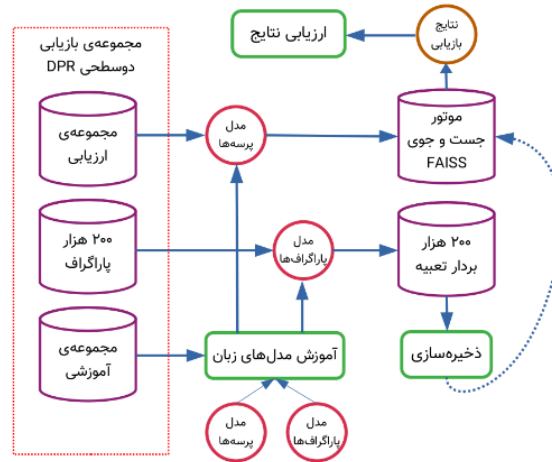
<sup>۱۳</sup> Stochastic Gradient Descent

<sup>۱۴</sup> Dense Passage Retrieval

استفاده از نمونه‌های منفی برای یک سؤال در نظر گرفته است. نمونه‌های منفی می‌توانند هرگونه سند تصادفی از مجموعه داده یا اسناد برتر بازگردانده شده توسط BM25 که حاوی پاسخ صحیح نیستند، باشد. از آنجایی که مدل DPR پایه تنها دو نوع رابطه «مرتبط» و «نامرتبط» را می‌شناسد، در ابتدا ساختار و نحوه‌ی آموزش مدل DPR بر روی مجموعه داده ارائه شده دوسطحی توضیح داده شده است. سپس در بخش بعدی تغییرات مورد نیاز در DPR برای درک و فهم مجموعه داده ارائه شده چهارسطحی مورد بحث قرار گرفته است.

آموزش و ارزیابی مدل DPR پایه در سه مرحله انجام می‌شود:

- آموزش دادن رمز گذارهای دو گانه با کمک مجموعه داده آموزشی
  - تعبیه کردن مجموعه‌ی اسناد با کمک رمز گذار مربوط به اسناد و ذخیره سازی مقادیر تعبیه شده با کمک ابزار FAISS (Johnson, Douze, and Jégou 2019)
  - بازیابی اطلاعات برای پرسه‌های داخل مجموعه‌ی ارزیابی با کمک رمز گذار مربوط به پرسه‌ها و FAISS و سپس ارزیابی کردن کارایی سامانه بر اساس نتایج حاصل شده
- در شکل ۲ می‌توان جایگاه این سه مرحله را مشاهده کرد. در این شکل می‌توان دید که علاوه بر مجموعه داده‌ی دوسطحی، دو نمونه از مدل زبان مورد نیاز می‌باشد تا بتوان آن‌ها را آموزش داد و سپس از آن‌ها برای تعبیه کردن اسناد و پرسه‌ها استفاده کرد.



شکل ۲: مراحل آموزش، تعبیه و ارزیابی در DPR

برای آموزش مدل بازیابی اطلاعات، یک مدل زبان فارسی و یا یک مدل زبان چندزبانه که از زبان فارسی نیز پشتیبانی کند مورد نیاز است. در این پژوهش، از مدل BERT چندزبانه<sup>۱۵</sup> استفاده شده است. همان‌طور که قبلاً نیز اشاره شد، جهت آموزش مدل، به دو نمونه از این مدل زبان نیاز می‌باشد. یک نمونه برای تعبیه کردن اسناد (پاراگراف) که آن را با  $E_p$  نشان می‌دهیم و نمونه‌ی دیگر برای تعبیه کردن پرسه‌ها که آن را با  $E_q$  نشان می‌دهیم. در زمان آموزش، یک دسته از داده‌های آموزشی انتخاب می‌شوند و همه به طور موازی پردازش می‌شوند. در هر دسته، هر پرسه شامل تعدادی سند مثبت و منفی می‌باشد که در مجموعه‌ی داده تعیین شده‌اند. علاوه بر اسناد منفی داخل مجموعه، اسناد مثبت پرسه‌های دیگر نیز می‌توانند به عنوان سند منفی پرسه جاری استفاده شوند (Karpukhin et al. 2020). پس از تعبیه کردن پرسه‌ها و اسناد، با استفاده از ضرب داخلی بردارها، شباهت میان هر پرسه و سند محاسبه می‌شود که آن را به صورت  $Sim(q, p)$  نشان می‌دهیم که  $q$  یک پرسه و  $p$  یک پاراگراف است و برای محاسبه‌ی آن از رابطه‌ی ۱ استفاده می‌شوند.

$$sim(q, p) = E_q(q)^T E_p(p) \quad (1)$$

<sup>۱۵</sup> BERT-base-multilingual-uncased

هرچه مقدار ضرب داخلی یک سند و پرسه بیشتر باشد، ارتباط میان آن دو قوی تر می باشد. برای نمونه، اگر فرض شود که برای آموزش مدل از دسته های سه تایی استفاده می شود و هر پرسه به همراه یک سند مثبت و یک سند منفی باشد، با کمک رابطه ی ۱ می توان ماتریس رسم شده در جدول ۱ را بدست آورد. در این جدول، سند  $d_1$  یک سند مثبت برای پرسه  $q_1$  می باشد؛ اما همین سند برای پرسه های  $q_2$  و  $q_3$  یک سند منفی محسوب می شود.

جدول ۱: ماتریس فرضی برای محاسبه ی شباهت سه پرسه

	سند مثبت برای پرسه ۱	سند منفی برای پرسه ۱	سند مثبت برای پرسه ۲	سند منفی برای پرسه ۲	سند مثبت برای پرسه ۳	سند منفی برای پرسه ۳
	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
$q_1$	$\text{sim}(q_1, d_1)$	$\text{sim}(q_1, d_2)$	$\text{sim}(q_1, d_3)$	$\text{sim}(q_1, d_4)$	$\text{sim}(q_1, d_5)$	$\text{sim}(q_1, d_6)$
$q_2$	$\text{sim}(q_2, d_1)$	$\text{sim}(q_2, d_2)$	$\text{sim}(q_2, d_3)$	$\text{sim}(q_2, d_4)$	$\text{sim}(q_2, d_5)$	$\text{sim}(q_2, d_6)$
$q_3$	$\text{sim}(q_3, d_1)$	$\text{sim}(q_3, d_2)$	$\text{sim}(q_3, d_3)$	$\text{sim}(q_3, d_4)$	$\text{sim}(q_3, d_5)$	$\text{sim}(q_3, d_6)$

سپس با کمک ماتریس محاسبه شده، به ازای هر سطر از ماتریس، تابع زیان محاسبه می شوند تا مقادیر زیان را به دست آوریم. رابطه ی ۲ این تابع زیان را نشان می دهد که Negative Log Likelihood (و یا به اختصار NLL) نام دارد (Karpukhin et al. 2020). در این رابطه،  $(q_i, p_i^+, p_{i-1}^-, \dots, p_{i-n}^-)$  نشان دهنده ی یک سطر از ماتریس با شاخص  $i$  است که شامل یک پرسه  $(q_i)$ ، یک سند مثبت  $(p_i^+)$  و چند سند منفی  $(p_{i-1}^-, \dots, p_{i-n}^-)$  می باشد. بخش کسری این رابطه همواره در بازه ی  $[0, 1]$  است و هرچه مقدار این کسر به ۱ نزدیک تر باشد، به این معنا است که مدل پیش بینی دقیق تری انجام داده و شباهت پرسه و سند مثبت را بالاتر ارزیابی کرده است. به طور مشابه، اگر مدل ما شباهت میان یک یا چند سند منفی با پرسه را بالا محاسبه و ارزیابی کند، مقدار مخرج این کسر افزایش پیدا می کند و مقدار نهایی کسر به صفر نزدیک تر می شود. مقدار شباهت محاسبه شده توسط بخش کسری این رابطه، با کمک عملیات لگاریتمی و منفی سازی آن، به یک مقدار زیان تبدیل می شود. شایان ذکر است که خروجی این رابطه همواره در بازه ی  $[0, \infty)$  می باشد.

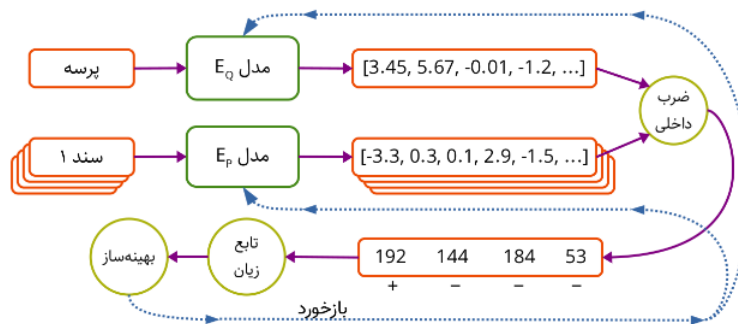
$$L(q_i, p_i^+, p_{i-1}^-, \dots, p_{i-n}^-) = -\text{Log} \frac{e^{\text{sim}(q_i, p_i^+)}}{e^{\text{sim}(q_i, p_i^+)} + \sum_{j=1}^n e^{\text{sim}(q_i, p_{i,j}^-)}} \quad (2)$$

با در دست داشتن مقادیر زیان، می توان به رمز گذار  $E_P$  (مربوط به پاراگراف ها) و همینطور رمز گذار  $E_Q$  (مربوط به پرسه ها) یک بازخورد از نحوه ی عملکرد آن ها داد. اگر مقدار تابع زیان نزدیک به صفر



باشد، پارامترهای مدل کمتر دستخوش تغییرات می‌شوند و هرچه مقدار تابع زیان بیشتر شود، تغییرات در پارامترهای مدل بیشتر خواهد شد. اعمال این تغییرات تا زمانی که مدل قادر به پیش‌بینی صحیح شباهت میان پرسه‌ها و اسناد مثبت شود ادامه می‌یابد. با کمک این تابع زیان می‌توان مدل را تنظیم دقیق کرد و در طی این مرحله، مدل یاد می‌گیرد تا برای اسناد و پرسه‌هایی که از منظر بازیابی اطلاعات مرتبط هستند، بردارهای مشابه تولید کند تا میزان شباهت بردارهای آن‌ها، بیشینه شود.

در شکل ۳ می‌توان فرآیند کلی آموزش (تنظیم دقیق) را مشاهده کرد. در این شکل (که در آن فرآیند آموزش ساده‌سازی شده است) می‌توان دید که یک پرسه، یک سند مثبت (سند ۱) و سه سند منفی (اسناد ۲ تا ۴) به رمز گذارهای  $E_p$  و  $E_q$  داده شده‌اند و به چندین بردار عددی تبدیل شده‌اند. سپس با کمک تابع شباهت ضرب داخلی، شباهت هر یک از بردارهای اسناد با بردار پرسه محاسبه شده و یک لیست از این شباهت‌ها تولید می‌شود که عنصر اول آن مربوط به سند مثبت و بقیه‌ی عناصر مربوط به اسناد منفی است. تابع زیان با دانستن جایگاه سند مثبت می‌تواند عملکرد رمز گذارها را ارزیابی کند. پس از آن با کمک بهینه‌ساز بازخوردی برای تنظیم پارامترهای دو مدل تولید شده و بدین ترتیب، مدل تنظیم دقیق می‌شود.



شکل ۳: فرآیند آموزش DPR

پس از آن که دو مدل زبان مورد استفاده در مرحله‌ی قبل آموزش دیدند، تمامی اسناد، تعبیه و شاخص‌گذاری می‌شوند. برای تعبیه کردن، لازم است تا از مدل مربوط به اسناد  $E_p$  استفاده کرده و تمامی اسناد را به بردارهای معادل آن نگاشت کنیم. پس از تعبیه شدن اسناد، با کمک ابزار FAISS باید تمام بردارها ذخیره شوند تا بتوان بر روی آن‌ها جست‌وجو انجام داد. این ابزار قادر است تا پس از

ذخیره‌سازی اسناد، یک بردار دریافت کند و بردارهای شبیه به آن بردار را با سرعت بسیار بالا بازیابی کند. در این مرحله می‌توان با کمک FAISS و مدل زبان  $E_Q$ ، بازیابی اطلاعات را انجام داد. برای ارزیابی سامانه، از بخش ارزیابی مجموعه‌داده که شامل سه تایی‌های (پرسه، پاسخ پرسه، سند مرتبط) می‌باشد استفاده می‌شود. با تعبیه کردن هر پرسه با کمک مدل  $E_Q$  یک بردار حاصل می‌شود و با وارد کردن این بردار به FAISS تعداد مشخصی سند که بردارهای مشابه دارند حاصل می‌شود. امتیاز نهایی این سندها همان ضرب داخلی دو بردار است و اسناد بر اساس این مقدار مرتب می‌شوند. حال می‌توان با کمک سنج‌های ارزیابی، کیفیت بازیابی اطلاعات این سامانه را سنجید. سنج‌های ارزیابی در بخش (۱-۵) معرفی شده‌اند.

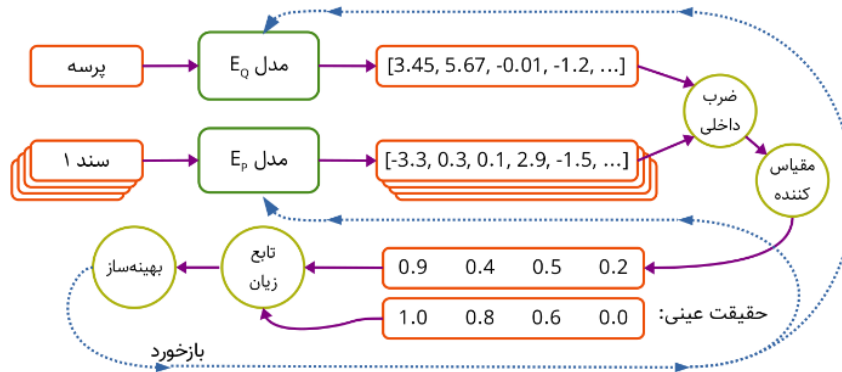
#### ۴-۲ مدل چندسطحی DPR

به منظور استفاده از مجموعه‌داده‌ی بازیابی اطلاعات چندسطحی، نیاز به تغییر دادن مدل DPR می‌باشد تا بتواند از سطوح تعریف شده در این مجموعه‌داده‌ی جدید استفاده کند. مهم‌ترین تغییری که بایستی در مدل DPR اعمال شود، بخش مربوط به تابع زیان است. این بخش باید قادر باشد تا به رتبه‌بندی مدل از اسناد چندسطحی، امتیاز دهد. در بخش (۴-۱) دیده شد که شباهت میان بردار اسناد و پرسه با کمک ضرب داخلی انجام می‌شود و تابع زیان مدل DPR پایه نیز، تابع NLL بود. در این بخش نیز از ضرب داخلی برای محاسبه‌ی شباهت استفاده شده است؛ اما تابع زبانی که در این بخش استفاده شده، باید قادر باشد تا یک بردار از شباهت‌ها (به صورت نرمال شده) دریافت کند و به این بردار امتیاز دهد. لذا برای مدل DPR چندسطحی، از تابع RankCosine (Xia et al. 2008) به عنوان تابع زیان استفاده شده است.

برای محاسبه‌ی تابع فوق، لازم است تا در ابتدا، شباهت‌های تولید شده توسط ضرب داخلی به بازه‌ی ۰ تا ۱ نگاشت شود. برای این منظور از یک مقیاس‌کننده‌ی min-max استفاده شده است. پس از آن با کمک تابع زیان، عملکرد سامانه سنجیده می‌شوند. در این بخش، تابع زیان علاوه بر دریافت بردار شباهت که حاصل محاسبات دو مدل زبان و ضرب داخلی بردارها می‌باشد، یک بردار حقیقت عینی<sup>۱۶</sup> نیز دریافت می‌کند که مقادیر آن به طور ضمنی توسط مجموعه‌داده تعیین می‌شوند و در واقع فاصله‌ی این دو بردار ورودی محاسبه می‌شود. بردار حقیقت با کمک میزان ارتباط اسناد با پرسه حاصل می‌شود. برای نمونه، اگر سند اول در بردار مشابهت، یک سند کاملاً مرتبط (سطح ۳) باشد، عنصر نظیر آن در بردار حقیقت

<sup>۱۶</sup> Ground Truth

۱ خواهد بود که معادل بالاترین مقدار ممکن برای امتیاز یک سند می‌باشد. به طور مشابه، عنصر نظیر یک سند نامرتب برابر ۰ خواهد بود که معادل پایین‌ترین مقدار ممکن برای امتیاز یک سند می‌باشد. در شکل ۴ می‌توان جایگاه بردار حقیقت در مدل بازیابی اطلاعات چندسطحی را مشاهده کرد. مقادیری که به ازای هر سطح در نظر گرفته می‌شوند، یک پارامتر برای مدل محسوب می‌شود و از همین جهت، برای یافتن مقدار مناسب، نیاز به تنظیم پارامتر می‌باشد.



شکل ۴: فرآیند آموزش DPR چندسطحی

برای محاسبه‌ی RankCosine ابتدا باید شباهت کسینوسی دو بردار ورودی محاسبه شود که این امر با کمک رابطه‌ی ۳ امکان پذیر می‌باشد. در این رابطه، بردار شباهت تولید شده توسط ضرب داخلی است و هر عنصر از این بردار، شباهت میان پرسه و یک سند را نشان می‌دهد. همان‌طور که قبلاً نیز اشاره شد، مقادیر این بردار به بازه‌ی ۰ الی ۱ مقیاس شده‌اند. بردار  $V_{Ground}$  همان بردار حقیقت است که هر عنصر آن، نشانگر سطح ارتباط سند با پرسه می‌باشد. مقادیر این بردار نیز در بازه‌ی ۰ الی ۱ می‌باشند. مقدار محاسبه شده توسط رابطه‌ی ۳ همواره میان ۱- و ۱ است.

$$\text{Cosin}(V_{Dot}, F_{Ground}) = \frac{V_{Dot}^T V_{Ground}}{\|V_{Dot}\| \|V_{Ground}\|} \quad (3)$$

برای محاسبه‌ی مقدار زیان از رابطه‌ی ۴ استفاده می‌شود (Xia et al. 2008). برای محاسبه‌ی این مقدار، ابتدا شباهت کسینوسی دو بردار محاسبه شده و سپس مقدار این شباهت به مقدار زیان تبدیل می‌شود.

$$L(V_{Dot}, V_{Ground}) = \frac{1}{2} [1 - \text{Cosine}(V_{Dot}, V_{Ground})] \quad (4)$$

با کمک این تابع زیان، می‌توان مدل را با کمک مجموعه داده‌ی چندسطحی آموزش داد. پس از فرآیند آموزش، نحوه‌ی استفاده از مدل‌ها دقیقاً مانند یک مدل دوسطحی می‌باشد. با کمک مدل‌های آموزش دیده، می‌توان اسناد و پرسه‌ها را تعیبه کرد، ضرب داخلی بردارهای حاصل شده را محاسبه نمود و اسناد را بر اساس امتیاز ضرب داخلی، مرتب کرد. در واقع در مرحله‌ی استفاده، معنا و مفهوم دوسطحی و چندسطحی بودن مدل‌ها از میان می‌رود و تنها رتبه‌بندی اسناد اهمیت خواهد داشت. به عبارت دیگر، تفاوت مدل چندسطحی و دوسطحی، تنها در بخش آموزش و مجموعه داده‌ی آموزشی می‌باشد و در مراحل تعیبه کردن اسناد و ارزیابی، از عملیات و مجموعه‌های ارزیابی کاملاً یکسانی استفاده می‌شود. علت این امر آن است که مدل‌ها (چه دوسطحی و چه چندسطحی) پس از مرحله‌ی آموزش، رفتار یکسانی از خود نشان می‌دهند و با دریافت یک پرسه و سند، یک عدد اعشاری به عنوان امتیاز تولید می‌کنند.

## ۵. ارزیابی

همانطور که پیشتر گفته شد، تفاوت مدل چندسطحی و دوسطحی، تنها در بخش آموزش و مجموعه داده‌ی آموزشی می‌باشد؛ و در مرحله‌ی ارزیابی، استفاده از مجموعه داده‌های دوسطحی و یا چندسطحی برای هر دو مدل امکان‌پذیر می‌باشد و این امر هیچ ارتباطی با نحوه‌ی آموزش دوسطحی یا چندسطحی مدل‌ها نخواهد داشت. برای نمونه، می‌توان یک مدل را به صورت دوسطحی آموزش داد و در ارزیابی آن، از یک مجموعه داده‌ی چندسطحی و یا دوسطحی استفاده کرد... در این پژوهش، برای داشتن یک ارزیابی صحیح و عادلانه، از یک مجموعه ارزیابی دوسطحی برای ارزیابی هر دو مدل استفاده شده است تا عملکرد مدل‌ها در شرایط برابر و یکسان سنجیده شود.

### ۱-۵ سنجه‌های ارزیابی مدل بازیابی اطلاعات

برای سنجش کیفیت بازیابی انجام شده در بازیابی اطلاعات، سنجه‌های مختلفی قابل استفاده هستند. از رایج‌ترین این سنجه‌ها می‌توان به سنجه‌ی Recall اشاره کرد که در رابطه‌ی ۵ نشان داده شده است (Mitra and Craswell 2018). در این رابطه  $(i, d) \in R_q$  نشان‌دهنده‌ی تمام اسناد بازیابی شده توسط سامانه‌ی بازیابی اطلاعات می‌باشد. مقدار  $d \in D$  نیز نشان‌دهنده‌ی هر سند مرتبط با پرسه  $q$  می‌باشد. مقادیر  $rel_q(d)$  همواره ۱ یا ۰ بوده و نشان‌دهنده‌ی ارتباط داشتن یا نداشتن سند مورد بررسی

$d$  با پرسه  $q$  می‌باشد. به طور خلاصه، این رابطه نسبت تعداد اسناد مثبت بازیابی شده (صورت کسر) را به تعداد کل اسناد مثبت داخل مجموعه داده (مخرج کسر) می‌سنجد.

$$Recall_q = \frac{\sum_{(i,d) \in R_q} rel_q(d)}{\sum_{d \in D} rel_q(d)} \quad (5)$$

از دیگر سنجه‌های مورد استفاده در بازیابی اطلاعات می‌توان به  $MRR^{17}$  اشاره کرد. برای محاسبه‌ی این سنجه، به ازای هر پرسه و نتایج بازیابی شده‌ی آن به کمک مدل، مقدار  $RR^{18}$  (که در رابطه‌ی ۶ نشان داده شده است) محاسبه می‌شود و سپس میانگین این مقادیر  $RR$  محاسبه و بعنوان  $MRR$  در نظر گرفته می‌شود (Mitra and Craswell 2018). در این رابطه،  $i$  را می‌توان به عنوان رتبه‌ی هر سند بازیابی شده تفسیر کرد. به عبارت دیگر، برای محاسبه‌ی  $RR$  هر سند، معکوس رتبه‌ی اولین سند مثبت محاسبه می‌شود و سپس با محاسبه‌ی میانگین برای این مقادیر  $RR$ ، مقدار  $MRR$  محاسبه می‌شود.

$$RR_q = \max_{(i,d) \in R_q} \frac{red_q(d)}{i} \quad (6)$$

سنجه‌ی دیگری که استفاده از آن در بازیابی اطلاعات رایج می‌باشد، سنجه‌ی Precision می‌باشد که بسیار شبیه به سنجه Recall بوده و در رابطه‌ی ۷ نشان داده شده است. در این رابطه  $|R_q|$  نشان‌دهنده‌ی تعداد کل اسناد بازیابی شده توسط سامانه می‌باشد. به عبارت دیگر، این سنجه نسبت تعداد اسناد مثبت بازیابی شده به تعداد کل اسناد بازیابی شده را می‌سنجد.

$$Precision_q = \frac{\sum_{(i,d) \in R_q} rel_q(d)}{|R_q|} \quad (7)$$

همچنین استفاده از گونه‌ی تغییر یافته‌ای از سنجه‌ی Precision که با نام Exact Match شناخته می‌شود در ارزیابی‌های سامانه‌های درک مطلب ماشینی بسیار رایج می‌باشد. در این سامانه‌ها، یافتن پاسخ دقیق یک پرسه در سند، یک موفقیت محسوب می‌شود و از همین رو، تمامی اسنادی که حاوی پاسخ باشند، مثبت فرض می‌شوند. به عبارت دیگر، برای محاسبه‌ی سنجه‌ی Exact Match کافی است نسبت اسناد بازیابی شده که حاوی پاسخ هستند به کل اسناد بازیابی شده محاسبه شود.

<sup>17</sup> Mean Reciprocal Rank

<sup>18</sup> Reciprocal Rank

در این پژوهش، برای ارزیابی از سنجه‌های Recall و MRR استفاده شده است. همینطور در صورتی که پاسخ هر پرسه در دسترس باشد، مقادیر Exact Match نیز محاسبه شده‌اند. لازم به ذکر است که تعداد اسناد بازیابی شده به ازای هر سنجه به صورت @num و در کنار نام سنجه نشان داده می‌شود. برای نمونه، اگر برای محاسبه‌ی MRR از ۱۰ سند و برای محاسبه‌ی Recall از ۱۰۰ سند استفاده شده باشد، آنها را به ترتیب با MRR@10 و Recall@100 نمایش می‌دهند. همینطور شایان ذکر است که مقادیر حاصل شده از هر یک از این سنجه‌ها یک مقدار اعشاری بین ۰ و ۱ می‌باشد. در این پژوهش، مقادیر حاصل شده از این سنجه‌ها در ۱۰۰ ضرب شده‌اند تا تفسیر کردن آن‌ها راحت‌تر و برحسب درصد باشد.

#### ۲-۵ نتایج آموزش و ارزیابی بر روی Trec-DL

مجموعه داده‌ی Trec-DL یک مجموعه داده‌ی چندسطحی انگلیسی است. در این پژوهش با کمک مجموعه داده‌های منتشر شده در سال‌های ۲۰۱۹ و ۲۰۲۰، دو مدل آموزش داده می‌شوند. مدل اول، یک مدل DPR دوسطحی است که تنها قادر به استفاده از سطوح ۰ و ۳ به عنوان اسناد مثبت و منفی است. مدل دوم، یک DPR چندسطحی است که می‌تواند بر خلاف مدل قبلی، علاوه بر سطوح ۰ و ۳، از سطوح ۱ و ۲ نیز استفاده کند. در ادامه‌ی این بخش، با کمک سنجه‌های Recall و MRR عملکرد این دو مدل مقایسه و بررسی می‌شوند. از آنجایی که پاسخ هر پرسه در مجموعه داده‌ی Trec-DL وجود ندارد، نمی‌توان سنجه‌هایی چون Exact Match را برای ارزیابی استفاده کرد.

برای ارزیابی دقیق‌تر، فرآیند آموزش هر یک از مدل‌های دوسطحی و چندسطحی سه مرتبه و با دانه‌های تصادفی<sup>۱۹</sup> مختلف انجام شده است.<sup>۲۰</sup> پس از آن میانگین این نتایج محاسبه و به عنوان نتیجه‌ی نهایی ارایه گردیده است. نتایج ارزیابی برای Recall و MRR بر روی مجموعه داده‌ی Trec-DL را می‌توان در جدول ۲ و جدول ۳ مشاهده نمود. همانطور که دیده می‌شود، مدل چندسطحی عملکرد بهتری نسبت به همتای دوسطحی خود دارد (۷,۰۴٪ برای Recall@100 و ۴,۰۲٪ برای MRR@100). این بهبود عملکرد نشان می‌دهد که تغییرات اعمال شده در مدل اصلی DPR و تبدیل آن به یک مدل چندسطحی، موثر واقع شده و مدل جدید چندسطحی DPR قادر بوده به درستی از اطلاعات اضافی

<sup>۱۹</sup> Random Seed

<sup>۲۰</sup> دانه تصادفی عددی تصادفی است که برای مقداردهی پارامترهای مدل و نوع درهم‌سازی داده‌های آموزشی و ... بکار میرود و از همین رو، تغییر دادن آن می‌تواند بر نحوه‌ی آموزش دیدن مدل و نتایج حاصل شده از مدل در مرحله‌ی ارزیابی موثر باشد. از آنجایی که آموزش دادن DPR یک فرآیند زمانبر می‌باشد، در این بخش برای ارزیابی، آموزش‌ها تنها با سه دانه‌ی تصادفی مختلف انجام شده.

فراهم شده (روابط سطح ۱ و ۲) در جهت افزایش دقت مدل بازیابی اطلاعات استفاده کند. حال که بر اساس نتایج بدست آمده، از عملکرد صحیح مدل چندسطحی اطمینان حاصل شد، در بخش بعدی به ارزیابی مجموعه داده‌ی فارسی و بررسی درستی مجموعه داده و نیز ارزیابی کارآیی روش‌های تولید آن پرداخته می‌شود.

جدول ۲: مقادیر Recall برای مجموعه داده‌های دوسطحی و چندسطحی Trec-DL

مجموعه داده	Recall@1	Recall@10	Recall@20	Recall@100
دوسطحی - اجرای اول	4.26	14.25	18.86	33.58
دوسطحی - اجرای دوم	4.95	16.37	21.68	37.58
دوسطحی - اجرای سوم	3.20	11.25	15.31	28.66
چهارسطحی - اجرای اول	3.88	13.73	18.31	31.22
چهارسطحی - اجرای دوم	8.28	22.38	27.72	43.33
چهارسطحی - اجرای سوم	9.22	24.98	31.23	46.38
دوسطحی - میانگین	4.14	13.96	18.62	33.27
چهارسطحی - میانگین	7.13	20.36	25.75	40.31

جدول ۳: مقادیر MRR برای مجموعه داده‌های دوسطحی و چندسطحی Trec-DL

مجموعه داده	MRR@1	MRR@10	MRR@20	MRR@100
دوسطحی - اجرای اول	4.26	6.97	7.29	7.64
دوسطحی - اجرای دوم	4.95	8.04	8.40	8.78
دوسطحی - اجرای سوم	3.20	5.37	5.65	5.95
چهارسطحی - اجرای اول	3.88	6.55	6.87	7.17
چهارسطحی - اجرای دوم	8.28	12.17	12.55	12.91
چهارسطحی - اجرای سوم	9.22	13.58	14.01	14.37
دوسطحی - میانگین	4.14	6.79	7.11	7.46
چهارسطحی - میانگین	7.13	10.77	11.14	11.48

### ۳-۵ نتایج آموزش و ارزیابی بر روی مجموعه داده‌ی فارسی

در بخش (۲-۳) روش‌های مربوط به ساخت یک مجموعه داده‌ی چندسطحی جدید به کمک قوانین دست‌ساز و طبقه‌بندی خودکار مورد بحث قرار گرفت. در حالت طبقه‌بندی خودکار، به منظور انتخاب بهترین طبقه‌بند جهت ساخت مجموعه داده چندسطحی جدید، تعدادی طبقه‌بند مرسوم آزموده شدند که دقت طبقه‌بندی آن‌ها در جدول ۴ ذکر شده است. همان‌طور که مشاهده می‌شود، از میان این طبقه‌بندها، جنگل تصادفی عملکرد بهتری را ارائه کرده است. به همین دلیل، در این پژوهش نیز از همین طبقه‌بند جهت ساخت مجموعه داده‌ی چندسطحی فارسی استفاده شده است.

جدول ۴: دقت طبقه‌بندهای آزموده شده برای طبقه بندی اسناد

F1 (macro)	F1 (micro)	طبقه بند
<b>92.88</b>	<b>91.10</b>	جنگل تصادفی
82.39	78.00	درخت تصمیم
59.89	50.60	ماشین بردار پشتیبان
44.75	41.20	SGD
28.34	29.60	نزديکترین همسایه
45.38	40.70	AdaBoost
58.43	48.60	شبکه عصبی (۴ لایه مخفی)

به منظور ارزیابی تأثیرات مثبت مجموعه داده چندسطحی در مقابل مجموعه داده دوسطحی در سامانه‌های بازیابی اطلاعات، عملکرد هر دو مجموعه داده چندسطحی فارسی ساخته شده (یکی با قوانین دست‌ساز و یکی با طبقه‌بند جنگل تصادفی) با عملکرد مجموعه داده دوسطحی فارسی (به عنوان مدل پایه) مقایسه شده است. در جدول ۵ می‌توان مقادیر Exact Match@100 برای مدل دوسطحی پایه، مدل چندسطحی آموزش دیده بر روی مجموعه داده‌ی مبتنی بر قوانین دست‌ساز (با عنوان «چندسطحی قانون-محور») و مدل چندسطحی آموزش دیده بر روی مجموعه داده‌ی ایجاد شده با کمک جنگل تصادفی (با عنوان «چندسطحی جنگل تصادفی») را مشاهده کرد. هر یک از این مقادیر، نشان‌دهنده‌ی میانگین هفت اجرای مختلف با دانه‌های متفاوت می‌باشند. همان‌طور که مشاهده می‌شود، هر دو مدل چندسطحی، عملکرد بهتری نسبت به همتای دوسطحی خود دارند. مدل چندسطحی قانون-محور بهبود ۱٫۸۷ درصدی Exact Match را حاصل می‌کند و یا به تعبیر دیگر، مدل چندسطحی توانسته به ازای هر پرسه، حدوداً ۲ سند بیشتر (که شامل پاسخ نیز می‌باشد) نسبت به مدل دوسطحی بازیابی کند. مشاهده‌ی



دیگری که در این جدول می‌توان داشت، برتری روش قانون-محور نسبت به روش جنگل تصادفی می‌باشد. در واقع علت این امر در استفاده از ویژگی‌های مرتبط با پاسخ در این روش می‌باشد.

جدول ۵: مقادیر Exact Match برای مجموعه داده‌های دوسطحی و چندسطحی فارسی

مجموعه داده	Exact Match @100
دوسطحی	6.21
چندسطحی قانون-محور	8.08
چندسطحی جنگل تصادفی	7.35

در نهایت در جدول ۶ و جدول ۷ می‌توان به ترتیب ارزیابی Recall و MRR برای سه مجموعه داده را مشاهده کرد. در این جداول می‌توان عملکرد بهتر مجموعه داده‌ها و مدل‌های چندسطحی (به جز Recall@100 برای مدل چندسطحی جنگل تصادفی) را مشاهده کرد. از آنجایی که در بخش (۵-۲) عملکرد صحیح مدل‌ها مورد بررسی قرار گرفت، می‌توان ادعا کرد که بهبود حاصل شده در این بخش مربوط به کارا بودن دو روش ارائه شده این پژوهش (روش مبتنی بر قانون و روش مبتنی بر طبقه‌بند، معرفی شده در بخش‌های ۱-۲-۳ و ۲-۳-۲) می‌باشد. به عبارت دیگر، با کمک این دو روش، می‌توان به طور موثر یک مجموعه داده بازیابی اطلاعات دوسطحی را به یک مجموعه‌ی چندسطحی تبدیل کرد.

جدول ۶: مقادیر Recall برای مجموعه داده‌های دوسطحی و چندسطحی فارسی

مجموعه داده	Recall@1	Recall@10	Recall@20	Recall@100
دوسطحی	14.62	39.73	48.04	67.81
چندسطحی قانون-محور	17.97	41.58	48.60	68.62
چندسطحی جنگل تصادفی	16.63	42.58	50.15	67.50

جدول ۷: مقادیر MRR برای مجموعه داده‌های دوسطحی و چندسطحی فارسی

مجموعه داده	MRR@1	MRR@10	MRR@20	MRR@100
دوسطحی	14.62	21.96	22.53	23.02
چندسطحی قانون-محور	17.97	24.74	25.23	25.71
چندسطحی جنگل تصادفی	16.63	24.36	24.89	25.33

## ۶. نتیجه گیری

در مباحث پردازش زبان طبیعی، زبان فارسی یک زبان کم منبع می باشد و نیازمندی های زیادی در آن وجود دارد. در چند سال اخیر و با پیشرفت مدل های زبان، سامانه های بازیابی اطلاعات زیادی توسط محققان بر پایه ی این ابزارهای جدید ارائه گردیده است؛ اما برای زبان فارسی، هیچ پژوهشی در این زمینه انجام نشده و در نتیجه، هیچ مجموعه داده مناسبی برای آن وجود ندارد.

در این پژوهش، پس از تبدیل یک مجموعه داده ی درک مطلب ماشینی به یک مجموعه داده ی بازیابی اطلاعات دوسطحی، دو روش مختلف برای تبدیل کردن مجموعه ی ساخته شده به یک مجموعه داده ی چندسطحی ارائه شد. سپس مدل DPR به عنوان مدل پایه ی این پژوهش مورد استفاده قرار گرفت و تغییرات لازم در آن برای درک مجموعه داده های چندسطحی مورد بحث قرار گرفت.

با کمک دو مجموعه داده ی چند سطحی جدید و مجموعه داده ی دوسطحی و همینطور مدل های دوسطحی و چندسطحی، می توان به ارزیابی راه کارهای ارائه شده پرداخت. ارزیابی های انجام شده بر روی این مجموعه ها نشان می دهد که در صورت فراهم بودن مجموعه داده ی چندسطحی، مدل بازیابی اطلاعات می تواند نتایج بهتری نسبت به مدل دوسطحی حاصل کند. همینطور با کمک روش های پیشنهاد شده در این پژوهش، می توان یک مجموعه ی دوسطحی را به یک مجموعه ی چندسطحی تبدیل کرد تا بتوان از مزایای مدل چندسطحی، بهره برد.

## ۷. منابع

- Abadani, Negin, Jamshid Mozafari, Afsaneh Fatemi, Mohammad Ali Nematbakhsh, and Arefeh Kazemi. 2021. "ParSQuAD: Machine Translated SQuAD Dataset for Persian Question Answering." in 2021 7th International Conference on Web Research (ICWR). IEEE.
- Ayoubi Sajjad & Davoodeh, Mohammad Yasin. 2021. "PersianQA: A Dataset for Persian Question Answering." GitHub Repository.
- Bajaj, Payal, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNameara, Bhaskar Mitra, and Tri Nguyen. 2016. "Ms Marco: A Human Generated Machine Reading Comprehension Dataset." ArXiv Preprint ArXiv:1611.09268.
- Craswell, Nick, Bhaskar Mitra, Emine Yilmaz, and Daniel Campos. 2021. "Overview of the TREC 2020 Deep Learning Track." ArXiv Preprint ArXiv:2102.07662.
- Craswell, Nick, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Jimmy Lin. 2021. "Ms Marco: Benchmarking Ranking Models in the Large-Data Regime." pp. 1566-76 in Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval.
- Craswell, Nick, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. "Overview of the TREC 2019 Deep Learning Track." ArXiv Preprint ArXiv:2003.07820.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. "BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding." pp. 4171-86 in Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). Minneapolis, Minnesota: Association for Computational Linguistics.

- Hashemi, Helia, Mohammad Aliannejadi, Hamed Zamani, and W. Bruce Croft. 2020. "ANTIQUE: A Non-Factoid Question Answering Benchmark." Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) 12036 LNCS:166–73. doi: 10.1007/978-3-030-45442-5\_21.
- Johnson, Jeff, Matthijs Douze, and Hervé Jégou. 2019. "Billion-Scale Similarity Search with GPUs." IEEE Transactions on Big Data 7(3):535–47.
- Karpukhin, Vladimir, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. "Dense Passage Retrieval for Open-Domain Question Answering." pp. 6769–81 in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP).
- Kazemi, Arefeh, Jamshid Mozafari, and Mohammad Ali Nematbakhsh. 2022. "PersianQuAD: The Native Question Answering Dataset for the Persian Language." IEEE Access 10:26045–57. doi: 10.1109/ACCESS.2022.3157289.
- Khashabi, Daniel, Arman Cohan, Siamak Shakeri, Pedram Hosseini, Pouya Pezeshkpour, Malihe Alikhani, Moin Aminnaseri, Marzieh Bitaab, Faeze Brahman, and Sarik Ghazarian. 2021. "ParsiNLU: A Suite of Language Understanding Challenges for Persian." Transactions of the Association for Computational Linguistics 9:1163–78.
- Kwiatkowski, Tom, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, and Kenton Lee. 2019. "Natural Questions: A Benchmark for Question Answering Research." Transactions of the Association for Computational Linguistics 7:453–66.
- Lin, Jimmy, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. "Pyserini: A Python Toolkit for Reproducible Information Retrieval Research with Sparse and Dense Representations." pp. 2356–62 in Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval.
- Liu, Ye, Kazuma Hashimoto, Yingbo Zhou, Semih Yavuz, Caiming Xiong, and S. Yu Philip. 2021. "Dense Hierarchical Retrieval for Open-Domain Question Answering." pp. 188–200 in Findings of the Association for Computational Linguistics: EMNLP 2021.
- Liu, Zhenghao, Kaitao Zhang, Chenyan Xiong, Zhiyuan Liu, and Maosong Sun. 2021. "OpenMatch: An Open Source Library for NEU-IR Research." pp. 2531–35 in Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval.
- Mitra, Bhaskar, and Nick Craswell. 2018. "An Introduction to Neural Information Retrieval." Foundations and Trends® in Information Retrieval 13(1):1–126.
- Qu, Yingqi, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2021. "RocketQA: An Optimized Training Approach to Dense Passage Retrieval for Open-Domain Question Answering." pp. 5835–47 in Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.
- Rajpurkar, Pranav, Robin Jia, and Percy Liang. 2018. "Know What You Don't Know: Unanswerable Questions for SQuAD." pp. 784–89 in Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers).
- Rajpurkar, Pranav, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. "SQuAD: 100,000+ Questions for Machine Comprehension of Text." pp. 2383–92 in Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing.
- Robertson, Stephen, and Hugo Zaragoza. 2009. "The Probabilistic Relevance Framework: BM25 and Beyond." Foundations and Trends® in Information Retrieval 3(4):333–89.
- Salton, Gerard, and Christopher Buckley. 1988. "Term-Weighting Approaches in Automatic Text Retrieval." Information Processing & Management 24(5):513–23.
- Xia, Fen, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. 2008. "Listwise Approach to Learning to Rank: Theory and Algorithm." pp. 1192–99 in Proceedings of the 25th international conference on Machine learning.



- Yang, Peilin, Hui Fang, and Jimmy Lin. 2017. "Anserini: Enabling the Use of Lucene for Information Retrieval Research." pp. 1253–56 in Proceedings of the 40th international ACM SIGIR conference on research and development in information retrieval.
- Zhang, Xinyu, Andrew Yates, and Jimmy Lin. 2020. "A Little Bit is Worse than None: Ranking with Limited Training Data." pp. 107–12 in Proceedings of SustainNLP: Workshop on Simple and Efficient Natural Language Processing.

## Multi-level Persian Dataset for Information Retrieval

Ali Abedzadeh

Master of Software Engineering, Faculty of Computer Engineering,  
University of Isfahan, Isfahan, Iran, [a.abedzadeh@eng.ui.ac.ir](mailto:a.abedzadeh@eng.ui.ac.ir)

Reza Ramezani

Assistant Professor, Faculty of Computer Engineering,  
University of Isfahan, Isfahan, Iran, [r.ramezani@eng.ui.ac.ir](mailto:r.ramezani@eng.ui.ac.ir)

Afsaneh Fatemi

Associate Professor, Faculty of Computer Engineering,  
University of Isfahan, Isfahan, Iran, [a\\_fatemi@eng.ui.ac.ir](mailto:a_fatemi@eng.ui.ac.ir)

**Abstract.** Information retrieval systems are an essential part of many smart systems. The applications of this research field include search engines such as Google and Bing, question-answering systems, modern databases, etc. An information retrieval system tries to retrieve documents related to a question/query. The retrieval is done from a large collection of documents, and the size of this collection can be from a few thousand documents to millions of documents. In recent years, a lot of research has been done to develop information retrieval systems using language models. However, in this research field, no research has been done for the Persian language. One of its main reasons is the lack of a suitable Persian dataset for training language models. In this research, first, a Persian dataset for information retrieval is presented. After that, methods for enriching this data set are investigated. This enrichment is done by defining multi-level relationships between a document and a question. In this regard, the new dataset can show the relationship between question and document in four levels (unrelated - related - highly related - completely related) instead of two levels (completely unrelated - completely related). The name of the generated dataset is PersianMLIR. Experiments show that by using multi-level relationships, the performance of the system improves for both Persian and English languages, where the improvement is 1.87% for the Persian language. The results conclude that enriching information retrieval datasets by increasing the number of relations between query and document lead to improving the performance of information retrieval systems.

**Keywords.** Information Retrieval, Language Models, Information Retrieval Dataset, Persian Dataset

### علی عابدزاده

متولد سال ۱۳۷۵ دارای مدرک تحصیلی کارشناسی ارشد در رشته مهندسی کامپیوتر از دانشگاه اصفهان است. عنوان پایان‌نامه ایشان «ساخت مجموعه داده چندسطحی جهت آموزش مدل‌های زبان برای زبان کم‌منبع فارسی» می‌باشد که تحت راهنمایی جناب آقای دکتر رضا رضانی و مشاوره سرکار خانم دکتر افسانه فاطمی به پایان رساندند.



### رضا رضانی

متولد سال ۱۳۶۸ دارای مدرک تحصیلی دکتری تخصصی در رشته مهندسی کامپیوتر گرایش نرم‌افزار از دانشگاه فردوسی مشهد است. ایشان هم‌اکنون استادیار گروه مهندسی نرم‌افزار دانشکده‌ی مهندسی کامپیوتر دانشگاه اصفهان است. پردازش زبان طبیعی، تحلیل داده و وب معنایی از جمله علایق پژوهشی وی است.



### افسانه فاطمی

متولد ۱۳۵۲ دارای مدرک تحصیلی دکتری تخصصی در رشته مهندسی کامپیوتر گرایش نرم‌افزار از دانشگاه اصفهان است. ایشان هم‌اکنون دانشیار گروه مهندسی نرم‌افزار دانشکده‌ی مهندسی کامپیوتر دانشگاه اصفهان است. سیستم‌های پیچیده، کلان‌داده، سیستم‌های پرسش‌پاسخ و سیستم‌های گفتگو از جمله علایق پژوهشی وی است.

