

سیستم هوشمند پرسش و پاسخ مبتنی بر ویکی‌پدیا در حوزه‌ی پزشکی

غزاله خرادپور

دانشجوی کارشناسی رشته علوم کامپیوتر، دانشکده علوم کامپیوتر دانشگاه صنعتی امیرکبیر، تهران، ایران

استاد راهنما: دکتر اکبری

چکیده: در این پروژه، سعی بر آن بود که با وجود سیستم پرسش و پاسخ مبتنی بر مقاله‌های ویکی‌پدیا^۱ غیرفارسی، داده‌هایی به زبان فارسی -چه به صورت اتوماتیک و چه به صورت دستی- جمع‌آوری شود و سپس به کمک سیستم‌های هوشمند موجود، سیستمی پیاده‌سازی کنیم که کاربران بتوانند با وارد کردن سؤالات پزشکی خود به زبان فارسی، پاسخی درخور همان سوال به همان زبان فارسی دریافت کنند.

کلمات کلیدی: سیستم پرسش و پاسخ، ویکی‌پدیا، پزشکی، دیتا، اتوماتیک، پردازش زبان طبیعی

۱ مقدمه

بزرگی برای کاربران باشد از آن جا که پزشک و اطلاعات پزشکی در اختیار تمام مردم نیست و حتی اگر در دسترس هم باشد، ممکن است در یافتن جواب‌های خود دچار گمراهی شوند. افراد با حل مشکلات پزشکی سطحی خود، امکانات محدود پزشکی را در اختیار افرادی قرار می‌دهند که به آن نیاز شدیدی دارند. از طرفی جواب دادن به سؤالات پزشکی نیازمند دانستن دانش بسیار است و هر فردی در شرایط عادی آن‌ها را نمی‌داند و وقت زیادی ندارد تا منابع مختلف را مطالعه کند، این سیستم می‌تواند بسیار کمک‌کننده باشد. علاوه بر این پزشکان نیز می‌توانند صحت تشخیص و دانسته‌های خود را از این طریق تثبیت کنند و دانش پزشکی و ارتباط با بیماران به سمت هوشمندتر شدن می‌رود. بنابراین سیستمی که بتواند به پرسش‌های کاربران پاسخ مناسب برگرداند، بسیار خدمت بزرگی است و برای آن‌ها کم‌هزینه است. تنها با وارد کردن سؤالات خود به زبان طبیعی، به جواب متناسب می‌رسند.

ما برای پیدا کردن جواب سؤالات پاسخ‌بلند خود از ویکی‌پدیا به عنوان منبع دایرةالمعارف‌گونه استفاده می‌کنیم. ویکی‌پدیا منبع قدرتمندی است که اطلاعات به‌روزی دارد و هر روز بهتر می‌شود؛ چون توسط عده زیادی از مردم نوشته و ویرایش می‌شود، کم‌تر حاوی اطلاعات غلط است و می‌توان سیستم‌های هوشمندی روی آن پیاده کرد.

پرسش و پاسخ یک زمینه تحقیقاتی در رشته علوم کامپیوتر و در حوزه بازیابی اطلاعات و پردازش زبان طبیعی^۲ است که هدف از آن طراحی سیستم‌هایی است که به‌طور خودکار به سؤالات مطرح شده توسط انسان در ساختار زبان طبیعی پاسخ تولید می‌کند. پیاده‌سازی سیستم پرسش و پاسخ، معمولاً به‌صورت یک برنامه رایانه‌ای است که پاسخ‌های خود را با پرس‌وجو از یک پایگاه داده ساخت یافته از دانش یا اطلاعات، معمولاً یک پایگاه دانش ایجاد می‌کند. به‌طور معمول، سیستم‌های پرسش و پاسخ می‌توانند پاسخ‌ها را از یک مجموعه بدون ساختار از اسناد زبان طبیعی دریافت کنند. از سری اسناد زبان طبیعی که برای دریافت جواب سؤالات از آن استفاده می‌شود، می‌توان به مجموعه‌ای محلی از متون مرجع اسناد و صفحات وب سازمان داخلی، گزارش‌های خبری خبرگزاری‌ها، مجموعه‌ای از صفحات ویکی‌پدیا و زیرمجموعه‌ای از صفحات وب جهانی اشاره کرد. برای دریافت جواب‌ها از اسناد دو نوع منبع دامنه محدود و دامنه باز وجود دارد که هر کدام شرایطی دارند و انتخاب ما اسناد دامنه محدود بوده است به این معنی که با یک دامنه خاصی از موضوعات سروکار دارند و موضوع مورد نظر ما حوزه پزشکی بوده است. در آینده آن را به دامنه باز تبدیل خواهیم کرد.

حال سیستم پرسش و پاسخ هوشمند در حوزه پزشکی می‌تواند کمک

categories	title	text	paragraph_text	link	answer	question
Infobox medical condition' زده: رده: اختلال	اسهال جرب	اسهال جرب (انگلیسی: Steatorrhea) یا جرب پرخال	اسهال به طور کلی، به عنوان افزایش ... وزن مدفوع (ب	https://fa.wikipedia.org/wiki/%D8%A7%D8%B3%D9%...	علت آن عدم جذب صحیح چربی‌ها در ... روده باریک است	دلیل ایجاد اسهال جرب چیست؟
Infobox medical condition' زده: رده: اختلال	اسهال جرب	اسهال جرب (انگلیسی: Steatorrhea) یا جرب پرخال	اسهال جرب با حجم زیاد و براق از ... عوارض اسهال‌ها	https://fa.wikipedia.org/wiki/%D8%A7%D8%B3%D9%...	اسهال جرب، سبب ایجاد اختلال شدید ... در فرایند جذب	عوارض اسهال جرب چیست؟
Infobox medical condition' زده: رده: بیماری	واریس مری	به (واریس مری (انگلیسی: Esophageal varices)، ووا	درمان با انشاع بالون یا تزریق واروسین ... وریدی	https://fa.wikipedia.org/wiki/%D9%88%D8%A7%D8%...	درمان با انشاع بالون یا تزریق واروسین وریدی ...	درمان واریس مری چگونه انجام می‌شود؟
Infobox medical condition' زده: رده: بیماری	واریس مری	به (واریس مری (انگلیسی: Esophageal varices)، ووا	به (واریس مری (انگلیسی: Esophageal varices)، ووا	https://fa.wikipedia.org/wiki/%D9%88%D8%A7%D8%...	خطر اصلی واریس مری، پارگی عروق ... متسع و خونریزی	خطر اصلی واریس مری چیست؟
های 2006 تا زده: بیماری	میزوپروستول	به (میزوپروستول (انگلیسی: Misoprostol)، سترکینی	سردرد، تهوع، کرامپ رحمی، دل ... درد، اسهال، نفخ و	https://fa.wikipedia.org/wiki/%D9%85%DB%8C%D8%...	سردرد، تهوع، کرامپ رحمی، دل ... درد، اسهال، نفخ و	عوارض جانبی میزوپروستول چیست؟
...

شکل ۱: نمونه‌ای از داده جمع‌آوری شده

۲ منابع داده

نحوه استخراج جواب از زمینه به ماشین استفاده می‌کنیم. داده‌ها دارای ساختاری به شکلی زیر هستند:

- سؤال: سؤالی بر اساس مطلب موجود در ویکی‌پدیا طرح شده است.
- جواب: جواب سؤال طرح‌شده به صورت عینی از قسمتی از پاراگراف‌های یکی از مطالب ویکی‌پدیا نوشته شده است.
- زمینه: پاراگرافی که از آن جواب استخراج شده، ذخیره شده است. این ستون از داده‌ها برای آموزش مدل بسیار حائز اهمیت است.
- تیترا: تیترا مطلبی از ویکی‌پدیا که سؤال و جواب از آن گرفته شده است.
- متن کل: متن تمام مقاله مورد نظر در ویکی‌پدیا ذخیره شده است.
- رده‌ها: رده‌هایی که در قسمت پایینی صفحه ویکی‌پدیا به صورت لیستی قرار داده شده‌اند، ذخیره شده است.

لازم به ذکر است که سه مورد انتهایی مجموعه داده‌ها به صورت اتوماتیک و با استفاده از کد استخراج شده است. نمونه‌ای داده‌ها در شکل ۱ آمده‌اند

۲-۳ مجموعه داده‌های ترجمه شده

از آن‌جا داده‌های جمع‌آوری شده به صورت دستی کافی نبودند، شروع به ترجمه مجموعه داده SQuAD [۱] کردیم و با این که تقریباً ربط زیادی به حوزه پزشکی نداشتند و API سیستم مترجم گوگل^۴ قادر به ترجمه تمام و کمال مجموعه داده‌ها نبود، این داده‌ها هم بلااستفاده ماندند و در مراحل بعدتر و برای بهتر کردن مدل آموزش داده‌شده به ماشین از آن‌ها استفاده خواهیم کرد. علاوه بر خود مجموعه داده‌های این مجموعه، ستون‌هایی به آن‌ها اضافه کردیم که شامل ترجمه‌هایی به زبان فارسی هستند و سوال، جواب، تیترا مقاله و زمینه به زبان فارسی ترجمه شده‌اند.

منابع اصلی ما داده‌های تمامی صفحات ویکی‌پدیا و داده‌های جمع‌آوری شده بودند که از داده‌های ویکی‌پدیا به عنوان منبعی برای استخراج جواب استفاده می‌کنیم و از داده‌های جمع‌آوری شده برای آموزش دادن مدلی که قادر است بهترین جواب برای سؤال پرسیده شده را برگرداند. چالش اصلی کار ما، پیدا کردن داده مناسب به زبان فارسی بود و از آن‌جا کار مشابهی قبلاً به زبان فارسی انجام نشده است و به طور کلی در پردازش زبان طبیعی به فارسی پیشینه خاصی وجود ندارد، داده مناسبی هم در اختیار نداشتیم. چندین منبع برای این نوع مجموعه داده‌ها وجود داشتند که هم محدود بودند و هم ساختاری به شکل آن چه ما می‌خواستیم نداشتند. کاری که ما می‌خواستیم انجام دهیم، چون پردازش زبان طبیعی بود و باید از گونه‌ای از شبکه عصبی بهره می‌گرفتیم، به دیتای بسیار زیادی برای آموزش نیاز داشتیم تا به هدفی که برای خود تعیین کرده بودیم، برسیم.

۲-۱ ویکی‌پدیا

از ویکی‌پدیا به عنوان منبع دانش استفاده می‌شود. حال برای استخراج تمامی مقالات ویکی‌پدیا باید حدود ۵ میلیون صفحه استخراج شود و تمامی شکل‌ها و ساختارهای مربوطه از بین برود و فقط متن به تنهایی در نظر گرفته و ذخیره شود.

۲-۲ مجموعه داده‌های جمع‌آوری شده

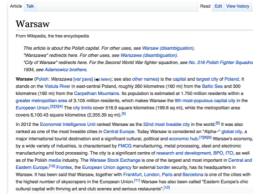
برای این که بتوانیم مدلی آموزش دهیم تا نتیجه نزدیک به مطلوب را به ما بدهد، نیاز به داده‌های بسیار زیادی از زبان فارسی داشتیم و از آن‌جا که مجموعه داده‌ها -خصوصاً در حوزه پزشکی- بسیار کم بود، شروع به جمع‌آوری هزار ردیف داده فارسی -با این که باز هم مقدارش کم بود- کردیم. این دیتاها به صورت دستی از سایت ویکی‌پدیا استخراج شدند و ستون‌های آن شامل سؤال، جواب، متن پاراگراف (زمینه^۳)، تیترا مقاله و متن کل مقاله است. از متن پاراگراف و کل مقاله برای یاد دادن

Open-domain QA SQuAD, TREC, WebQuestions, WikiMovies

Q: How many of Warsaw's inhabitants spoke Polish in 1933?

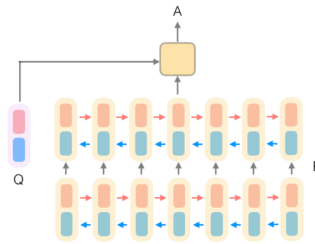


**Document
Retriever**



**Document
Reader**

833,500



شکل ۲: شیوه کار سیستم DrQA

پروژه در شکل ۲ نشان داده شده است.

۳ مروری بر کارهای پیشین

برای سیستم پرسش و پاسخ به زبان فارسی، هیچ کاری انجام نشده است چون همان طور که قبلاً اشاره کردیم، مشکل کمبود داده و پیچیدگی های زبان فارسی چالش های زیادی را ایجاد می کند. به زبان های دیگر خصوصاً انگلیسی کارهای زیادی انجام شده است.

کارهای مشابه دیگری به عنوان مثال توسط Ryu et al. [۲] انجام شده است که سیستم پرسش و پاسخ دامنه بازی طراحی کرده است که از ویکی پدیا به عنوان منبع مدل خود استفاده می کند و متن مقاله ای را با جواب های منطبق دیگر مانند جعبه های اطلاعات^۵، ساختار مقاله، ساختار دسته بندی و تعاریف ترکیب می کند. به طور مشابه Ahn et al. [۳] هم که از ویکی پدیا به عنوان منبع متن با ترکیب دیگر منابع استفاده می کنند. علاوه بر آن، Buscaldi and Rosso [۴] از ویکی پدیا به عنوان منبع دانش خود استفاده می کنند اما به جای آن که جواب ها را از روی آن بیابند، از آن برای ارزیابی کردن جواب تولید شده توسط خود سیستم استفاده می کنند و از دسته بندی های ویکی پدیا برای اطمینان از مجموعه الگوهایی که برای جواب پیش بینی شده، بهره می برد.

سیستم های پیشرفته تری از پرسش و پاسخ وجود دارند که از ویکی پدیا و وب استفاده می کنند؛ مانند QuASE^۵، Microsoft's AskMSR [۵]، IBM's DeepQA [۶] که هر کدام از روش های خاص خود برای پیاده سازی این سیستم استفاده کرده اند. که یکی از بزرگ ترین و مشابه ترین کارها به کار ما، پروژه ای به نام DrQA [۷] است و کاری را انجام داده است که برنامه داشتیم انجام دهیم. بدنه این پروژه به گونه ای است که با وارد کردن داده فارسی می توان مدلی را آموزش داد و از نتیجه آن استفاده کرد. شکل کلی این

۱-۳ DrQA چیست؟

ک سیستم درک مطلب از طریق خواندن است که از ویکی پدیا به عنوان تنها منبع خود برای یافتن جواب استفاده می کند به طوری که انگار شخصی از روی دایرةالمعارف جواب سؤالات را می دهد. البته این سیستم به گونه ای طراحی شده است که می توان هر منبعی را به عنوان منبع پیدا کردن جواب سؤالات برای آن انتخاب کرد و قابلیت مقیاس پذیری دارد. این سیستم پس از دریافت سؤال، به دنبال جوابی برای آن سؤال می گردد. تمرکز این پروژه بر روی سؤال و جواب های طولانی (نه یک کلمه ای) است و در این مسیر با چالش هایی هم چون پیدا کردن مقاله های مرتبط و پیدا کردن جواب از روی آن ها دست و پنجه نرم می کند. این پروژه از ویکی پدیای انگلیسی - از آن جا که بسیار غنی و با جزئیات است - به عنوان منبع جواب های خود استفاده می کند.

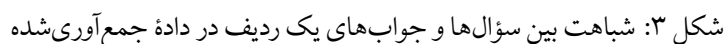
DrQA از دو قسمت اصلی تشکیل شده است که قسمت اول آن تمام مقالات ذخیره شده را می خواند و از بین آن ۵ مقاله ای پتانسیل دارا بودن جواب را در خود دارند، بر می گرداند و قسمت دوم آن همان هسته اصلی است که مدل روی آن آموزش دیده است.

این پروژه از شبکه عصبی مکرر چندلایه^۶ که در حالت پیش فرض روی مجموعه داده های SQuAD آموزش داده شده است، به عنوان پردازش کننده مطالب استفاده می کند و در نهایت جواب های منتخب را با امتیاز حساب شده توسط مدل ارائه می دهد و بهترین آن را به عنوان جواب بر می گرداند.

• Document Retriever: تمام داده ها و مقالات ویکی پدیا

^۶multi-layer recurrent neural network machine comprehension

^۵infoboxes

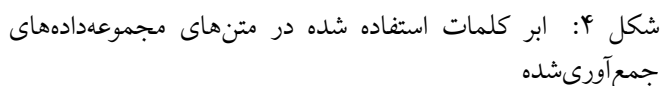


به صورت ماتریس های وزن داده شده TF-IDF ذخیره می شوند. پس از آن که سؤال دریافت شد، توسط ادغام TF-IDF و n-grams - حتی سریع تر و با دقت بهتر از موتور جستجوی خود وبکم، پدیا- مقالات مرتبط را پیدا می کند.

- ## ۴ پیش‌پردازش‌های انجام‌شده روی داده‌ها

بخش اعظم این پروژه جمع‌آوری داده‌ها بوده است که برای درک بهتر ساختارها و اجزای آن یک سری پیش‌پردازش‌ها انجام داده‌ایم تا بتوانیم با چه جنس داده و با چه توزیعی از کلمات مواجه هستیم. یکی از کارهای اصلی‌ای که برای پیش‌پردازش داده‌ها در فرآیند پردازش زبان طبیعی انجام می‌شود، حذف کلمات کلیدی^۸ است که انجام دادن آن در این پروژه بی‌معنی است. تنها کلمات غیرفارسی را از آن حذف کردیم و الگوهای زبانی را در آن بهبود بخشیدیم. دو نمونه از آمارهایی که در مورد مجموعه داده‌ها رسم کرده‌ایم، در شکل ۴ و ۳ آمده‌اند.

برای این که ارتباط بین سؤالات و جواب‌های هر ردیف از داده‌ها را بدانیم، با استفاده از روش TF-IDF متن سؤال‌ها و جواب‌ها به بردار تبدیل کردیم و سپس با روش Cosine Similarity شباهت بین سؤال‌ها و جواب‌های هر ردیف را بررسی کردیم. سؤال‌ها و جواب‌ها تا حدی



۵ نتیجه

این پروژه با بیشتر شدن داده‌ها، به امتیاز بیشتری خواهد رسید و هدف ما این است که این پروژه را در زمینه‌های دیگری به غیر از پزشکی نیز گسترش دهیم و در مراحل بعدتر به الگوریتم‌های موجود در این زمینه بهبود بخشیم.

- [1] K. L. P. L. Pranav Rajpurkar, Jian Zhang, "SQuAD: 100,000+ Questions for Machine Comprehension of Text," 2017.
- [2] H.-K. K. Pum-Mo Ryu, Myung-Gil Jang, "Open domain question answering using Wikipedia-based knowledge model," 2014.
- [3] G. M. K. M. M. d. R. S. S. David Ahn, Valentin Jijkoun, "Using Wikipedia at the TREC QA Track.," 2004.
- [4] P. R. Davide Buscaldi, "Mining knowledge from Wikipedia for the question answering task," 2006.
- [5] M. B. Eric Brill, Susan Dumais, "An Analysis of the AskMSR Question-Answering System," 2002.
- [6] J. □. Petr Baudiš, "Modeling of the Question Answering Task in the YodaQA System," 2015.
- [7] J. W. A. B. Danqi Chen, Adam Fisch, "Reading Wikipedia to Answer Open-Domain Questions," 2017.