

Solving Persian Crosswords with Natural Language Processing Techniques

¹Mohammadreza Pakzadian*, ²Mehrnoush Shamsfard

¹ Bachelor of Computer Engineering, Faculty of Computer Science and Engineering, Shahid Beheshti University, Tehran, Iran
pakzadianmrp@gmail.com

² Associate Professor, Faculty of Computer Science and Engineering, Shahid Beheshti University, Tehran, Iran
m-shams@sbu.ac.ir

Abstract

In this article, we introduce solutions for solving crossword puzzles by machine using natural language processing techniques. This task is divided into two subtasks of finding possible answers for each table description and then selecting the target word and placing it in the table. The first subtask, which is dedicated to finding the word from its description, has many other uses as in text generation and paraphrasing. For this purpose, we used a combination of different methods, including searching and finding semantic similarities on the data of previously solved tables, searching in dictionary and Wikipedia articles, using a masked language model, and finding related words in Farsnet and the Farsiyar tool. The results show that the combination of these methods has a better result (82% recall) compared to their individual implementation.

In the next subtask, we give the list of possible answers to a constraint-satisfaction search algorithm to choose the correct answer that can be placed in the table, taking into account the constraints of the table, and fill the empty cells in the best way and solve the crossword. The overall evaluation shows 80.22% precision and 68.86% recall in solving the crossword puzzle.

Keywords: crossword solving, constraint satisfaction problem, natural language processing, description understanding

حل جدول کلمات متقاطع فارسی با تکنیک‌های پردازش زبان طبیعی

محمد رضا پاکزادیان^{۱*}، مهرنوش شمس فرد^۲

^۱ کارشناسی مهندسی کامپیوتر، دانشکده مهندسی و علوم کامپیوتر، دانشگاه شهید بهشتی، تهران
pakzadianmrp@gmail.com

^۲ دانشیار، گروه هوش مصنوعی، رباتیک و رایانش شناختی، دانشکده مهندسی و علوم کامپیوتر، دانشگاه شهید بهشتی، تهران
m-shams@sbu.ac.ir

چکیده

در این مقاله، به معرفی راهکارهایی برای حل جدول کلمات متقاطع توسط ماشین با استفاده از تکنیک‌های پردازش زبان طبیعی می‌پردازیم. این کار به دو زیروظیفه یافتن کلمات محتمل برای هر شرح جدول و سپس انتخاب کلمه هدف و جایگذاری در جدول تقسیم می‌شود. زیروظیفه اول که به یافتن کلمه از روی توصیف آن اختصاص دارد دارای کاربردهای متعدد دیگر نیز هست. به این منظور از ترکیبی از روش‌های مختلف شامل جستجو و شباهت‌یابی معنایی روی دادگان جداول حل شده قبلی، جستجو در فرهنگ لغات و دانشنامه ویکی‌پدیا، استفاده از مدل زبانی ماسک شده و یافتن کلمات مرتبط در فارسی‌نت و ابزار فارسی‌یار استفاده نمودیم. نتایج نشان می‌دهند ترکیب این روش‌ها نسبت به اجرای تک تک آن‌ها نتیجه بهتری (فراخوانی ۸۲٪) در برداشته است. در زیروظیفه بعد لیست پاسخ‌های محتمل را به یک الگوریتم جستجو با ارضاء قیود می‌دهیم تا با در نظر گرفتن قیود جدول از میان پاسخ‌ها، پاسخ درستی که می‌تواند در جدول قرار گیرد را انتخاب کرده و خانه‌های خالی را به بهترین شکل پر کند و جدول را حل کند. ارزیابی کل نشانگر دقت ۸۰.۲۲٪ و فراخوانی ۶۸.۸۶٪ در حل جدول کلمات متقاطع است.

کلمات کلیدی

حل جدول کلمات متقاطع، جستجو با ارضاء قیود، پردازش زبان طبیعی، فهم تعریف.

پاسخ شرح‌ها استفاده می‌کند اما Proverb و WebCrow از تعدادی پیمانانه^۱ سفارشی برای این منظور استفاده می‌کنند. WebCrow همچنین از یک موتور جستجوی وب نیز برای این کار استفاده می‌کند. همچنین Berkely از یک مدل پرسش و پاسخ^۲ برای یافتن پاسخ‌ها استفاده می‌کند. برای بخش حل جدول، Proverb و WebCrow از الگوریتم انتشار باورهای حلقه‌ای^۳ در ترکیب با الگوریتم جستجوی *A استفاده می‌کنند در حالی که Dr.fill از یک الگوریتم جستجوی عمق اول^۴ اصلاح شده استفاده می‌کند و پس از آن با استفاده از یک الگوریتم جستجوی محلی پاسخ‌های قرار گرفته در جدول را اصلاح می‌کند و با استفاده از یک تابع مکاشفه^۵ به این جدول‌های اصلاح شده یک امتیاز نسبت می‌دهد. همچنین Berkely با استفاده از الگوریتم انتشار باورهای حلقه‌ای به همراه جستجوی محلی به پر کردن جدول می‌پردازد. در سال ۲۰۲۱ برای اولین بار در مسابقات جدول کلمات متقاطع آمریکایی، ترکیبی از سیستم‌های Berkely و Dr.fill توانست از تمام رقبای انسانی خود بهتر عمل کند.

۱- مقدمه

حل کردن جدول کلمات همواره به عنوان یک سرگرمی جذاب برای افراد شناخته می‌شود و بسیاری از افراد در اوقات فراغت خود به حل جدول می‌پردازند. همچنین پیشرفت‌های فناوری در زمینه هوش مصنوعی نشان داده است که کامپیوترها می‌توانند در بازی‌هایی مانند شطرنج، انسان‌ها را به چالش بکشند و حتی برنده شوند. از این رو حل کردن جدول کلمات متقاطع توسط کامپیوترها می‌تواند یک مسئله جذاب و چالش برانگیز برای افرادی که در حوزه علوم کامپیوتری کار می‌کنند، باشد. لازمه حل جدول کلمات متقاطع، درک معنای شرح‌های داده شده و یافتن کلمه‌ای با محدودیت‌های مشخص طول و حروف معلوم شده، متناظر با شرح است.

از جمله نرم افزارهای مشابهی که در زمینه حل جدول کلمات وجود دارند می‌توان به [3] WebCrow، [4] Dr.fill، [5] Proverb و Berkeley [6] Crossword Solver اشاره کرد که همگی برای حل جدول به زبان انگلیسی طراحی شده‌اند. Dr.fill از یک روش شبیه به TFIDF برای یافتن

۷	۶	۵	۴	۳	۲	۱	
ن	گ	ل		ر	ب	ج	۱
ی	ر	ی	ش		د	و	۲
ت		ت	م	ا	ه	ش	۳
ر	ا	ی		م	ی	د	۴
و	ج	م	ز	ر		ا	۵
ژ	ر		ن	و	ر	د	۶
ن	ت	م		د	ز	ن	۷

۲- طرح مسأله

مسأله‌ای که ما به دنبال حل آن هستیم این است که راهکارهایی برای حل جدول کلمات متقاطع فارسی توسط ماشین ارائه کنیم. برای این منظور ابتدا باید ساختاری برای تعریف جدول مشخص کنیم و سپس با به کارگیری تکنیک‌ها و روش‌های هوش مصنوعی و پردازش زبان‌های طبیعی پاسخ مورد نظر برای هریک از شرح‌های جدول را پیدا کرده و با انتخاب کلمه هدف و جایگذاری آن، جدول را کامل کنیم. نمونه‌ای از جدول کلمات متقاطع مورد نظر در شکل (۱) قابل مشاهده است.

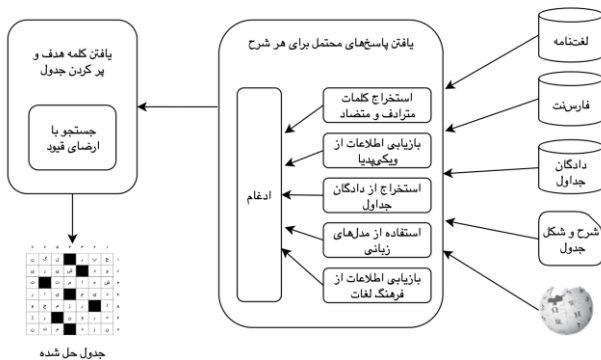
۳- راه حل پیشنهادی

حل جدول شامل دو مرحله اصلی است. (۱) یافتن کلمات کاندیدا به عنوان پاسخ محتمل با توجه به هر شرح (۲) انتخاب کلمه صحیح از میان کاندیداها با توجه به جدول و کلمات دیگر موجود در آن. شمای کلی مراحل مورد استفاده در حل جدول در شکل (۲) قابل مشاهده است.

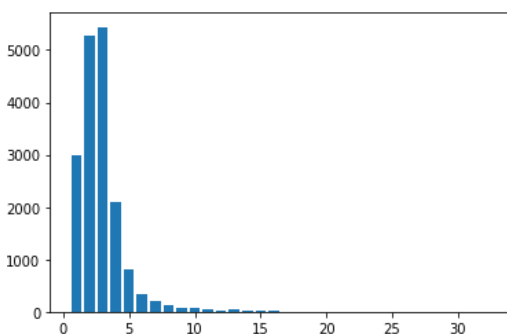
۱-۳- تهیه دادگان مورد نیاز

دادگان مورد نیاز را به دو دسته تقسیم می‌کنیم. دسته اول دادگان آموزش و آزمون برای جدول کلمات متقاطع است. به این منظور، با نوشتن یک برنامه خزشگر بر روی سایت جدول‌باب[۱]، تعداد ۱۰۸۸ عدد جدول از این سایت دریافت و ذخیره شدند. تهیه این مجموعه داده از دو جهت حائز اهمیت است. اول اینکه از بخشی از این داده‌ها می‌توان به عنوان مجموعه دادگان آزمون برای بررسی عملکرد نرم‌افزار استفاده کرد و دوم اینکه در بسیاری از جدول‌ها سوالات تکراری وجود دارند و از این مجموعه داده می‌توان به عنوان داده‌های آموزش برای یافتن پاسخ اینگونه سوالات استفاده کرد. جداول استخراج شده به صورت تصادفی به دو مجموعه ۱۰۰۰ تایی برای آموزش و ۸۸ تایی برای آزمون تقسیم شدند. نمودارهای فراوانی شرح‌ها برحسب تعداد کلمات موجود در آن و فراوانی پاسخ‌ها برحسب طول کلمه در این مجموعه دادگان در شکل‌های (۳) و (۴) آمده است. در این مجموعه بلندترین شرح حاوی ۳۲ کلمه و بلندترین کلمه پاسخ به طول ۱۶ است.

- عمودی:
 ۱- پیوند دادن
 ۲- قرض و دین - برنج
 ۳- گلایی
 ۴- حس بویایی - یار مرد
 ۵- عنصر شیمیایی فلزی
 ۶- مخفف اگر - مزد
 ۷- ازت
- افقی:
 ۱- زور - یار آفتابه
 ۲- دوستی - رنگ سفید مات
 ۳- شجاعت و دلیری
 ۴- کشت بارانی - دوست
 ۵- جنگاور
 ۶- داخل - از لوازم آرایش
 ۷- پیش و جلو - درون حاشیه
- شکل (۱): نمونه‌ای از یک جدول کلمات متقاطع فارسی



شکل (۲): مراحل مورد استفاده در حل جدول



شکل (۳): نمودار فراوانی شرح‌ها برحسب تعداد کلمات

در ادامه ابتدا به دادگان مورد استفاده پرداخته و سپس روش به کار گرفته شده برای هریک از دو مرحله فوق را شرح می‌دهیم. دسته دوم دادگان پایگاه‌های داده و دانش در مورد کلمات، توصیف آن‌ها و مثال‌هایی از کاربردشان است. جمع آوری اطلاعات از لغت‌نامه دهخدا[۲]، ویکی‌پدیای فارسی، فارس‌نت[7] و فارسی‌یار[8] در این دسته قرار دارند. در این راستا با خزش بر روی سایت لغت‌نامه دهخدا تعداد ۳۴۳۳۱۸ کلمه و معنی متناظر با آن جمع آوری شد. از فارس‌نت (وردنت[9] فارسی) نیز ۱۰۰ هزار مدخل واژگانی با توضیح، مثال و روابط (ترادف، تضاد، شمول و ...) آن‌ها با کلمات دیگر استخراج شد. همچنین از مجموعه دادگان ویکی‌پدیا فارسی که حاوی بیش از یک میلیون مقاله فارسی است نیز استفاده شد. علاوه بر این دادگان، خواهیم دید که از ابزار استخراج رابطه معنایی کلمات در جعبه ابزار فارسی‌یار نیز بهره گرفته شده است.

در این جستجو پرس و جوی^۱ مورد استفاده همان شرح جدول است. جستجو بر روی کل اسناد انجام میگردد و ۵- بهترین نتیجه جستجو بعنوان اسناد کاندیدا بازگردانده می شود. سپس از میان این اسناد در ویژگی های title ، RedirectList و links به دنبال کلماتی با طول مورد نظر می گردیم و تمامی این کلمات را در مجموعه پاسخ های احتمالی قرار می دهیم. برای محاسبه احتمال درستی هر پاسخ، از فرمول (۲) استفاده می کنیم.

$$p = \frac{score}{100} * c \quad (2)$$

که در آن score امتیازی است که ابزار الستیک سرچ به هر کدام از این اسناد نسبت می دهد و c ضریب اهمیت ویژگی مورد نظر در آن سند است و مقدار آن برای ویژگی title برابر ۰.۹۹ و برای سایر ویژگی ها برابر ۰.۹ است. این روش برای یافتن پاسخ پرسش هایی مشابه پرسش های زیر عملکرد خوبی از خود نشان می دهد:

- تنظیم کننده قند خون (انسولین)
- از شهرهای استان فارس (استهبان)
- عنصر شیمیایی فلزی (لیتیم)

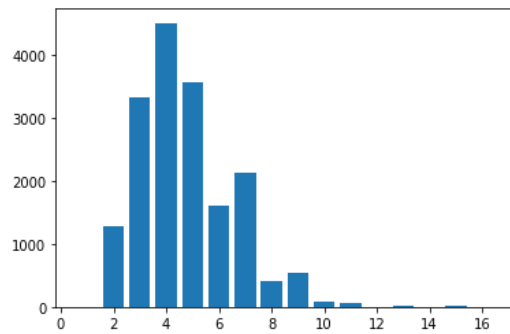
۳-۲-۳- استخراج از دادگان جداول

در بسیاری از مواقع سوال های جدول تکراری هستند و افرادی که جدول های زیادی حل می کنند با دانستن این موضوع می توانند عملکرد بهتری در حل جدول داشته باشند. وظیفه این پیمانانه این است که با استفاده از مجموعه داده ای که از جدول های سایت جدول یاب تهیه شد پاسخ اینگونه سوال های تکراری را پیدا کند و به حل سریع تر جدول کمک کند. این مجموعه داده حاوی ۱۷۶۶۶ ردیف داده متمایز که از ۱۰۰۰ عدد جدول استخراج شده اند می باشد. این روش با جستجو بر روی این مجموعه داده، شرح های مشابه با شرح مورد نظر را پیدا می کند و پاسخ آن ها را به عنوان پاسخ احتمالی برای این شرح در نظر می گیرد. پیاده سازی این پیمانانه به دو روش انجام شد که در ادامه به توضیح هر کدام از روش ها می پردازیم.

یکی از روش هایی که می توان برای پیاده سازی این پیمانانه استفاده کرد، استفاده از ابزارهای جستجوی متنی مانند الستیک سرچ برای یافتن شرح های مشابه است. برای این کار ابتدا ویژگی های سند مربوطه را در این ابزار تعریف کردیم و سپس داده های مورد نظر را وارد این ابزار کردیم.

سپس یک پرس و جو بر روی این اسناد انجام می دهیم و از میان اسنادی که در پاسخ این پرس و جو داده می شوند ۱۰- بهترین نتیجه جستجو را انتخاب می کنیم و پاسخ آن ها را به عنوان لیست پاسخ های احتمالی در نظر می گیریم. سپس این پاسخ ها را برحسب طول کلمه ای که به دنبال آن هستیم، محدود می کنیم و برای محاسبه اینکه هر کدام از این پاسخ ها با چه احتمالی می توانند پاسخ مورد نظر ما باشند، از فرمول (۳) استفاده می کنیم.

$$p_i = \frac{score_i}{\sum score} \quad (3)$$



شکل (۴): نمودار فراوانی پاسخ ها برحسب طول کلمه

۳-۲-۲- یافتن پاسخ های محتمل برای هر شرح

برای یافتن مجموعه پاسخ های محتمل برای هر شرح جدول از ترکیبی از چند روش استفاده کردیم. هر کدام از این روش ها علاوه بر پیدا کردن مجموعه پاسخ ها، به هر عضو این مجموعه یک عدد بین صفر و یک نسبت می دهند که نشان دهنده این است که این پاسخ با چه احتمالی می تواند پاسخ مورد نظر ما باشد. در ادامه به توضیح هر کدام از این روش ها که در قالب یک پیمانانه در سیستم پیاده سازی شده اند، می پردازیم.

۳-۲-۱- استخراج کلمات مترادف و متضاد

این روش که مناسب یافتن پاسخ برای شرح هایی با یک کلمه یا یک عبارت ساده است، با استفاده از فارسن ت و فارسی یار کلمات مترادف یا متضاد با کلمه یا عبارت شرح را پیدا می کند و پس از اعمال قید طول کلمه، مجموعه ای از پاسخ های محتمل را برمی گرداند. الگوی شرح های مورد پوشش به شکل زیر است:

- <کلمه> مثال: دوستی (ود)، مزد (اجرت)
- <کلمه> و <کلمه> مثال: شجاعت و دلیری (شهامت)
- مخالف <کلمه> مثال: مخالف کوچک (بزرگ)
- متضاد <کلمه> مثال: متضاد سریع (کند)

احتمال هر کدام از این پاسخ ها از طریق فرمول (۱) به دست می آید که در آن L برابر با تعداد کلمات موجود در مجموعه پاسخ های محتمل است.

$$p = \frac{1}{L} \quad (1)$$

۳-۲-۲- بازیابی اطلاعات از ویکی پدیا

بسیاری از پرسش های جدول به گونه ای هستند که پاسخدهی به آنها نیاز به دانش در حوزه خاصی مانند علوم، جغرافیا، تاریخ و... دارد. پاسخ اینگونه پرسش ها معمولاً با جستجو در دانشنامه هایی مانند ویکی پدیا قابل استخراج است. در این بخش، از جستجو روی مجموعه دادگان ویکی پدیای فارسی با استفاده از ابزار الستیک سرچ^۲ که بر پایه موتور جستجوی لوسین [10] توسعه یافته است، برای یافتن پاسخ اینگونه پرسش ها استفاده نمودیم.

امتیازی که مدل زبانی به هر کدام از نماینده‌ها نسبت می‌دهد، برای هر کدام از پاسخ‌ها احتمال را محاسبه می‌کنیم.

۳-۲-۵- بازیابی اطلاعات از فرهنگ لغات

وظیفه این بخش، بازیابی اطلاعات از فرهنگ لغات است. در این پیمانده اطلاعات لغت‌نامه دهخدا در ابزار الستیک سرچ وارد شد تا با استفاده از این ابزار بتوانیم بر روی معنی کلمات در لغت‌نامه دهخدا جستجوی متنی انجام دهیم و پاسخ مورد انتظار را پیدا کنیم.

در این جستجو پرس و جوی مورد استفاده همان شرح جدول است. جستجو بر روی کل اسناد انجام میگردد و ۱۰- بهترین نتیجه جستجو بعنوان اسناد کاندیدا بازگردانده می‌شود و کلمه‌ای که سند در مورد آن است به شرطی که قید طول کلمه را ارضا کند، به عنوان پاسخ احتمالی در نظر گرفته می‌شود. همچنین از فرمول ۳ برای محاسبه احتمال هر کدام از پاسخ‌ها استفاده می‌کنیم.

۳-۲-۶- ادغام پاسخ‌های پیشنهادی

پس از آنکه هر کدام از پیمانده‌ها لیست پاسخ‌های پیشنهادی خود برای هر کدام از شرح‌های جدول را ارائه کردند، این لیست‌ها با یکدیگر ادغام می‌شوند و مقدار احتمالی که برای هر پاسخ ارائه شده بود نیز برابر با بیشینه مقدار احتمال ارائه شده برای این پاسخ در پیمانده‌ها می‌شود.

۳-۳- پر کردن پاسخ‌ها در جدول

در این مرحله برای پر کردن پاسخ‌های احتمالی در جدول و پیدا کردن بهترین حالت پر کردن جدول، از الگوریتم جستجو با ارضای قیود^{۱۳} استفاده کردیم. در این الگوریتم از یک تابع ارزیابی برای مقایسه و ارزیابی حالت‌های مختلف جدول‌های پر شده و یک تابع مکاشفه‌ای برای پیش بینی امتیاز حالت نهایی جدول، برای جدول‌هایی که در حین جستجو به دست می‌آیند استفاده کردیم. فرمول (۴) تابع ارزیابی و فرمول (۵) تابع مکاشفه‌ای را نشان می‌دهد.

$$f(x) = a * \frac{N_{fq}}{N_q} + b * \frac{N_{cb}}{N_b} + c * \frac{\sum p_{x_i}(v_i)}{N_{fq}} \quad (4)$$

که در آن N_{fq} تعداد سوالات پر شده در جدول، N_{cb} تعداد کل سوالات جدول، N_{cb} تعداد خانه‌هایی از جدول که محل تقاطع دو سوال بوده‌اند و هر دو سوال پاسخ داده شده‌اند، N_b تعداد کل خانه‌های غیر سیاه جدول و $\sum p_{x_i}(v_i)$ مجموع احتمال ارائه شده برای پاسخ‌هایی که در جدول پر شده‌اند، می‌باشد.

$$h(x) = a * \frac{N_{fl}}{N_q} + b * \frac{N'_{cb}}{N_b} + c * \frac{\sum \max(p_{x_i}(v_i))}{N_{fq}} \quad (5)$$

که در آن N_{fl} تعداد سوالات پاسخ داده نشده که مجموعه پاسخ‌های محتمل آن‌ها شامل حداقل یک عضو است، N'_{cb} تعداد خانه‌هایی از جدول

که در آن score امتیازی است که ابزار الستیک سرچ به هر کدام از این اسناد نسبت می‌دهد.

روش دیگری که با استفاده از آن می‌توان این جستجو را انجام داد، استفاده از مدل‌های زبانی^۹ برای پیدا کردن جملات مشابه با استفاده از جستجوی معنایی^{۱۰} است. برای این منظور می‌توان از مدل [SBERT[11]] که بر پایه مدل زبانی [BERT[12]] و برای یافتن شباهت جملات ساخته شده است، استفاده کرد. مزیت این روش نسبت به روش قبل این است که در روش قبل معنای کلمات در جستجو نقشی ندارند و فقط شکل نوشتاری آن‌ها در نظر گرفته می‌شود اما در مدل‌های زبانی علاوه بر شکل نوشتاری، معنای کلمات نیز در جستجو نقش دارند.

با توجه به اینکه ساخت و آموزش یک مدل زبانی نیاز به مقدار زیادی داده و صرف منابع و زمان زیادی دارد، برای این کار از مدل‌های زبانی از پیش آموزش داده شده استفاده کردیم. به همین منظور مدل [ParsBERT[13]] را روی مجموعه دادگان جداول ریزتنظیم^{۱۱} کردیم. سپس با استفاده از مشابهت یابی معنایی، شرح‌های مشابه با شرح مورد نظر را پیدا می‌کنیم. از میان تمام شرح‌های مشابه، ۱۰ عدد از بهترین شرح‌هایی که شباهت بالای ۷۵ درصد با شرح مورد نظر ما دارند و همچنین قید طول پاسخ را ارضا می‌کنند را انتخاب می‌کنیم تا پاسخ آن‌ها به عنوان پاسخ محتمل در نظر گرفته شود. برای نسبت دادن احتمال به هر کدام از این پاسخ‌ها، از همان فرمول ۳ استفاده می‌کنیم و مقدار score در این فرمول برابر میزان شباهت در نظر گرفته می‌شود.

۳-۲-۴- استفاده از مدل‌های زبانی

یکی دیگر از روش‌هایی که با استفاده از آن می‌توان پاسخ یک شرح را به دست آورد استفاده از تکنیک‌های پرامپتینگ و مدل‌های زبانی بافتاری است. در این روش با استفاده از ۱۰ الگوی از پیش تعریف شده، از روی شرح جدول یک جمله می‌سازیم و پاسخ مورد نظر را به شکل جای خالی در شرح قرار می‌دهیم و از این مدل زبانی می‌خواهیم تا با توجه به معنای جمله جای خالی را با کلمه مناسب پر کند. چند نمونه از الگوها و جملات ساخته شده با استفاده از آن‌ها به شکل زیر است:

- <پاسخ> <شرح> است.

مثال: [MASK] سازمان آموزشی، علمی و فرهنگی سازمان ملل متحد است. (یونسکو)

- <شرح> <پاسخ> است.

مثال: آنچه مورد لزوم و احتیاج باشد، [MASK] است. (ضروری)

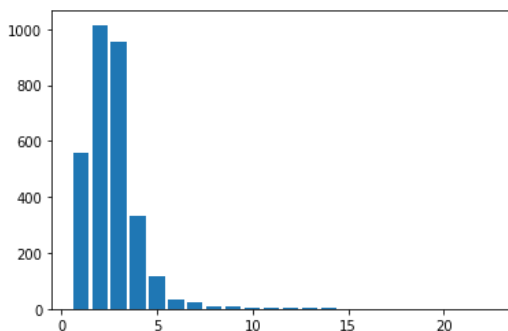
- <پاسخ> معادل <شرح> است.

مثال: مروت معادل [MASK] است. (جوانمردی)

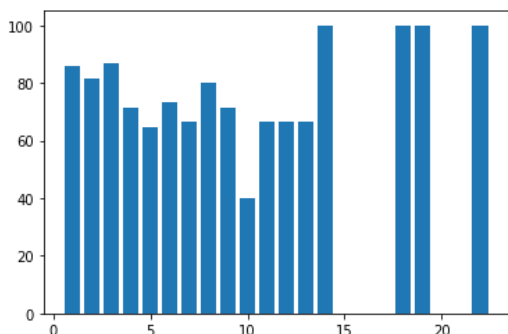
با توجه به اینکه هزینه آموزش یک مدل زبانی از پایه بسیار زیاد است، در این بخش ما از مدل زبانی ParsBERT که از قبل بر روی حجم زیادی از متون فارسی آموزش داده شده است، استفاده کردیم و این مدل را بر اساس مجموعه داده‌ای که داشتیم ریزتنظیم کردیم. سپس ۲۰ نماینده برتر برای پر کردن جای خالی را که قید طول قید طول کلمه مورد نظر را ارضا می‌کنند، به عنوان مجموعه پاسخ‌های محتمل انتخاب می‌کنیم و با استفاده از فرمول ۳ و

جدول (۱): نتیجه ارزیابی پیمانه‌ها

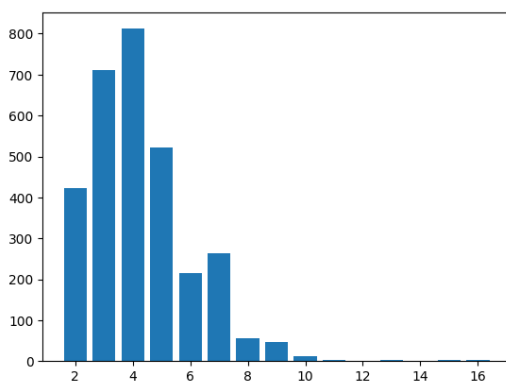
فراخوانی	نام پیمانه
٪۷۰.۸۱	استخراج کلمات مترادف و متضاد (فارسی‌نت)
٪۷۲.۵۳	استخراج کلمات مترادف و متضاد (فارسی‌یار)
٪۲۰.۹۰	بازیابی اطلاعات از ویکی‌پدیا
٪۷۳.۵۲	جستجوی الستیک روی دادگان جداول
٪۷۳.۴۰	شباهت‌یابی معنایی روی دادگان جداول
٪۱۰.۳۳	پرامپتینگ و مدل زبانی
٪۷.۷۱	بازیابی اطلاعات از فرهنگ لغات
٪۸۲.۱۹	تمام پیمانه‌ها در کنار هم



شکل (۵): نمودار فراوانی تعداد کلمات شرح در مجموعه دادگان آزمون



شکل (۶): نمودار فراوانی پیمانه‌ها بر حسب تعداد کلمات شرح



شکل (۷): نمودار فراوانی طول کلمه پاسخ در مجموعه دادگان آزمون

محل تقاطع دو سوال هستند و می‌توانند توسط سوال‌های پاسخ داده نشده پر شوند و $\sum \max(p(V_i))$ مجموع بیشینه احتمال ارائه شده در مجموعه پاسخ‌های محتمل برای سوالات پاسخ داده نشده، می‌باشد. در فرمول تابع ارزیابی و مکاشفه‌ای مقادیر a ، b و c ضرایبی هستند که میزان اهمیت هر بخش از تابع را مشخص می‌کنند و مقادیر در نظر گرفته شده برای آن‌ها که با آزمایش‌های تجربی بدست آمده به ترتیب ۱، ۴، و ۲ است.

در الگوریتم جستجو، سازگاری گره از قبل توسط پیمانه‌هایی که پاسخ‌های محتمل را تولید می‌کردند چک می‌شود. به این صورت که این پیمانه‌ها بر اساس قید طول کلمات که توسط جدول مشخص شده است، فقط کلمات با طول مناسب را نگه می‌دارند و سایر کلمات را از دامنه پاسخ‌های محتمل حذف می‌کنند. در مرحله انتخاب گره برای گسترش، در هر مرحله جدولی که مقدار تابع ارزیابی بیشتری داشته باشد برای گسترش و پر کردن سایر سوالات انتخاب می‌شود. سپس از میان سوالاتی که هنوز پر نشده‌اند، سوالی که بیشترین تعداد خانه‌های مشترک با سایرین را دارد انتخاب می‌شود تا پاسخ آن در جدول قرار گیرد و با توجه به اینکه یکی از سوالات در جدول قرار گرفته و قیود جدیدی در جدول اضافه شده‌اند، مجموعه پاسخ‌های احتمالی سایر سوالات مطابق این قیود محدود می‌شوند.

برای محدودتر کردن فضای جستجو، حلت‌های نامطلوب را هرس می‌کنیم. به این صورت که در هنگام جستجو همیشه بهترین جدولی که تا کنون یافته‌ایم یعنی جدولی که بیشترین مقدار تابع ارزیابی را داشته، نگه می‌داریم و اگر جدول‌هایی که در هنگام گسترش ایجاد می‌شوند، مقدار تابع مکاشفه‌ای آن‌ها کمتر از بهترین حلت باشد، این حلت‌ها را دور ریخته و هرس می‌شوند.

۳-۴- ارزیابی عملکرد پیمانه‌ها

برای ارزیابی عملکرد پیمانه‌ها از ۸۸ جدولی که به عنوان دادگان آزمون کنار گذاشته بودیم، استفاده می‌کنیم که این مجموعه داده حاوی ۸۱۰ ردیف داده متمایز است. برای هر کدام از پیمانه‌ها معیار فراخوانی^{۱۲} را مطابق فرمول (۶) محاسبه می‌کنیم.

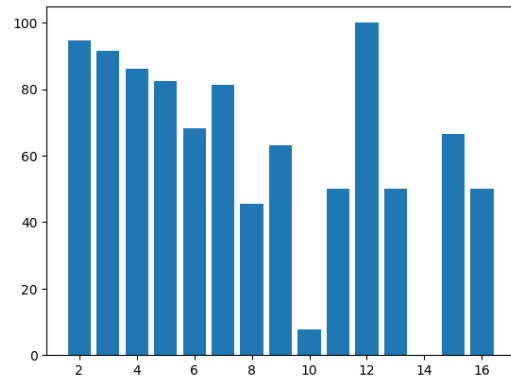
$$(۶) \quad \text{تعداد سوال‌هایی که پاسخ صحیح در لیست پاسخ‌های احتمالی وجود داشت} \\ \text{تعداد کل سوالات} = \text{فراخوانی}$$

نتیجه ارزیابی پیمانه‌ها در جدول (۱) قابل مشاهده است. در شکل (۵) نمودار فراوانی تعداد کلمات شرح‌ها در مجموعه دادگان آزمون و در شکل (۶) نمودار فراوانی تمام پیمانه‌ها در کنار هم بر حسب تعداد کلمات موجود در شرح قابل مشاهده است.

همچنین در شکل (۷) نمودار فراوانی طول پاسخ در مجموعه دادگان آزمون و در شکل (۸) نمودار فراوانی پیمانه‌ها بر حسب طول پاسخ قابل مشاهده است. همان طور که در این نمودارها مشخص است با افزایش طول پاسخ، فراخوانی کاهش می‌یابد زیرا پاسخ‌هایی که طول بیشتری دارند معمولاً مرکب از دو یا چند کلمه هستند و پیدا کردن آن‌ها سخت‌تر است.

7	6	5	4	3	2	1
ج	ب	ا	ر	ل	ی	ن
و	ا	د	ش	ا	ز	ی
ش	ا	ه	ا	م	ت	ا
ی	ا	م	ا	ا	ا	ر
و	ا	ج	م	ز	ا	و
د	ر	و	ن	ا	ر	ژ
ن	ز	ا	د	ا	ن	

شکل (۹): نمونه‌ای از عملکرد نرم افزار در حل جدول



شکل (۸): نمودار فراخوانی پیمانه‌ها برحسب طول کلمه پاسخ

۳-۵- ارزیابی عملکرد نرم افزار

برای ارزیابی عملکرد کلی نرم افزار حل جدول، از همان مجموعه داده آزمون که در بخش قبل توضیح داده شد استفاده کردیم و دو معیار فراخوانی و دقت ۱۴ را مطابق فرمول‌های (۷) و (۸) محاسبه کردیم.

$$\text{فراخوانی} = \frac{\text{تعداد خانه‌هایی از جدول که به درستی مقداردهی شده‌اند}}{\text{تعداد کل خانه‌های جدول}} \quad (۶)$$

$$\text{دقت} = \frac{\text{تعداد خانه‌هایی از جدول که به درستی مقداردهی شده‌اند}}{\text{تعداد خانه‌های پر شده}} \quad (۷)$$

پس از ارزیابی‌های انجام شده مشخص شد که نرم افزار ما به صورت میانگین، فراخوانی ۶۸.۸۶ درصد و دقت ۸۰.۲۲ درصد در حل جدول دارد. به عنوان نمونه عملکرد نرم افزار بر روی جدول شکل (۱) در شکل (۹) قابل مشاهده است.

۴- نتیجه گیری

حل کردن جدول کلمات متقاطع به کمک هوش مصنوعی مسئله‌ای چالش برانگیز و پیچیده است که نیازمند اطلاعات و داده در زمینه‌های مختلف است. برای اینکه بتوانیم یک نرم افزار حل کننده جدول کلمات با دقت خیلی بالا داشته باشیم، نیازمند مجموعه داده‌ای در ابعاد یک موتور جستجو هستیم تا بتوانیم در زمینه‌های مختلف اطلاعات داشته باشیم و همچنین باید بتوانیم به شکل مناسبی از این مجموعه داده اطلاعات مورد نظر را بازیابی کنیم که داشتن این حجم از داده و بازیابی مناسب آن نیازمند زمان، منابع و هزینه بسیار زیادی است.

به همین دلیل در این پروژه ما بر روی بخش‌های کوچک‌تر و هدفمندتری از داده متمرکز شدیم تا بتوانیم نیازمندی خود را برطرف کنیم. به همین منظور در این پروژه از مجموعه دادگان لغت‌نامه دهخدا، مجموعه مقالات فارسی ویکی‌پدیا، مجموعه جدول‌های سایت جدول‌یاب، فارس‌نت و فارسی‌ار استفاده کردیم تا بتوانیم تا حد قابل قبولی منبع دانش مورد نیاز برای حل جدول را شکل دهیم و سپس با استفاده از تکنیک‌های بازیابی اطلاعات و مدل‌های زبانی به دریافت اطلاعات از این مجموعه دانش پرداختیم. سپس با استفاده از الگوریتم

جستجو با ارضای قیود نسبت به حل جدول و پیدا کردن پاسخ مناسب برای هر کدام از سوالات اقدام کردیم. و در نهایت به فراخوانی ۶۸.۸۶ درصد و دقت ۸۰.۲۲ درصد دست پیدا کردیم.

برای کارهای آتی می‌توان بهبودهایی در بخش‌های مختلف نرم افزار ایجاد کرد تا عملکرد بهتری داشته باشد. یکی از این بهبودها استفاده از مجموعه دادگان وسیع‌تر و بزرگ‌تر است. مورد دیگری که قابل بهبود است، مدل‌های زبانی استفاده شده هستند که می‌توان با افزایش داده‌ها^{۱۵} و استفاده ترکیبی از مجموعه داده‌ها عملکرد آن‌ها را بهتر کرد. همچنین می‌توان به مرحله ادغام پاسخ‌ها یک مرحله یادگیری اضافه کرد تا بتوان لیست پاسخ‌های محتمل نهایی را به شکل مناسبی رتبه بندی کرد [14]. علاوه بر این موارد می‌توان بهبودهایی در الگوریتم جستجو، مانند یافتن توابع ارزیابی و مکاشفه‌های مناسب‌تر و همچنین هرس بهتر فضای جستجو، ارائه کرد.

مراجع

- [۱] جدول‌یاب، حل جدول آنلاین، مرداد ۱۳۹۶، <https://www.jadvalyab.ir>
- [۲] دهخدا، علی‌اکبر، لغت‌نامه دهخدا، انتشارات دانشگاه تهران، ۱۳۷۷، (نسخه برخط در آدرس اینترنتی <https://dehkhoda.ut.ac.ir>)
- [3] Ernandes, Marco, Angelini, Giovanni, Gori, Marco, "WebCrow: a web-based system for crossword solving.", AAAI'05: Proceedings of the 20th national conference on Artificial intelligence, Vol. 3, 2005
- [4] Ginsberg, Matthew L., "Dr. Fill: Crosswords and an implemented solver for singly weighted CSPs", Journal of Artificial Intelligence Research, Vol. 42, Issue 1, 2011
- [5] Littman, Michael L., Keim, Greg A., Shazeer, Noam, "A probabilistic approach to solving crossword puzzles." Artificial Intelligence, Vol. 134, Issue 1-2, 2002
- [6] Wallace E., Tomlin N., Xu A., Yang K., Pathak E., Ginsberg M., Klein D., "Automated Crossword Solving", Annual Meeting of the Association for Computational Linguistics, 2022
- [7] Shamsfard, M., Hesabi, A., Fadaei, H., Mansoori, N., Famian, A., Bagherbeigi, S., Fekri, E., Monshizadeh, M., Assi, S. M., "Semi automatic development of farsnet: the persian wordnet", Proceedings of 5th global WordNet conference, Mumbai, India, 2010, vol. 29
- [8] Khezry, Behrouz, Asgarian, Ehsan, "FarsiYar: Persian Text mining and text processing tools", 2018, <https://text-mining.ir>

- [9] Fellbaum Ch., "WordNet: An Electronic Lexical Database.", MIT Press, 1998
- [10] McCandless M., Hatcher E., Gospod-netic O., "Lucene in Action, Second Edition: Covers Apache Lucene 3.0", Manning Publications Co., Greenwich, CT, USA, 2010
- [11] Reimers, Nils, Gurevych, Iryna, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks", Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing, 2019
- [12] Devlin J., Chang M., Lee K., Toutanova K., "BERT: Pre-training of deep bidirectional transformers for language understanding." In NAACL, 2019
- [13] Farahani, M., Gharachorloo, M., Farahani, M., Manthouri, M., "ParsBERT: Transformer-based Model for Persian Language Understanding", Neural Processing Letters, Vol. 53, Issue 6, 2021
- [14] Barlacchi G., Nicosia M., Moschitti A., "Learning to Rank Answer Candidates for Automatic Resolution of Crossword Puzzles", Proceedings of the Eighteenth Conference on Computational Natural Language Learning, 2014.

زیر نویس

-
- ¹module
²question answering model
³loopy belief propagation
⁴Depth First Search (DFS)
⁵heuristic
⁶crawler
⁷elasticsearch
⁸query
⁹language model
¹⁰semantic search
¹¹finetune
¹²Constraint Satisfaction Problem (CSP)
¹³recall
¹⁴precision
¹⁵Data augmentation