

تخمین کیفیت پاسخ با استفاده از یادگیری جمعی در جوامع پرسش و پاسخ

معصومه خالقیان<sup>۱\*</sup> و میرمحسن پدram<sup>۲</sup>

<sup>۱</sup> دانشجوی مقطع کارشناسی ارشد رشته هوش مصنوعی از دانشگاه خوارزمی، masoumehkhaleghian@khu.ac.ir

<sup>۲</sup> دانشیار گروه مهندسی برق و کامپیوتر دانشکده فنی برق خوارزمی، pedram@khu.ac.ir

چکیده: جوامع پرسش و پاسخ به بخشی جدایی ناپذیر از زندگی روزمره انسان‌ها تبدیل شده‌اند. با وجود موفقیت این پلتفرم‌ها، بسیاری از پرسش‌ها بدون انتخاب پاسخ بهینه باقی می‌مانند، زیرا پرسش‌گران هیچ پاسخی را به عنوان بهترین گزینه تأیید نمی‌کنند. این وضعیت منجر به اتلاف زمان کاربران در جستجوی پاسخ‌های مناسب می‌شود. با رشد مداوم این جوامع، شناسایی دقیق بهترین پاسخ‌ها به چالشی اساسی تبدیل شده است. در این پژوهش، یک مدل یادگیری جمعی وزن دار ارائه شده است که چهار یادگیرنده پایه قدرتمند شامل *RoBERTa*، *XLNet*، جنگل‌های تصادفی و رأی‌گیری اکثریتی را ترکیب می‌کند. فرایادگیر این مدل که مبتنی بر درخت تصمیم است، با وزن‌دهی مناسب به خروجی یادگیرنده‌های پایه، دقت مدل را به سطح بی‌نظیر ۹۹.۰۳٪ در مجموعه داده آزمون ارتقا داده است. نتایج نشان می‌دهند مدل پیشنهادی می‌تواند نقش مهمی در بهبود کیفیت پاسخ‌ها، افزایش اعتماد به اطلاعات، و ارتقای تجربه کاربری در سیستم‌های پرسش و پاسخ ایفا کند.

کلیدواژه‌ها: تشخیص بهترین پاسخ، یادگیری جمعی، مدل‌های زبانی، پردازش زبان طبیعی، سیستم‌های پرسش و پاسخ

## ۱- مقدمه

### RoBERTa برای پردازش ویژگی‌های متنی و

مدل‌هایی همچون جنگل‌های تصادفی برای تحلیل ویژگی‌های غیرمتنی مورد استفاده قرار گرفته‌اند. همچنین، ترکیب این مدل‌ها از طریق روش‌های پشته‌ای و استفاده از فرایادگیر مبتنی بر درخت تصمیم، امکان دستیابی به دقتی فراتر از توانایی مدل‌های منفرد را فراهم کرده است.

این تحقیق تلاش دارد با ارائه مدل‌های ترکیبی پیشرفته و بهبود دقت پیش‌بینی‌ها، چالش‌های موجود در جوامع پرسش و پاسخ آنلاین را کاهش داده و به ارتقای کیفیت اطلاعات، افزایش رضایت کاربران و بهبود تجربه کاربری کمک کند.

## ۲- بیان مسأله و اهداف پژوهش

پلتفرم‌های پرسش و پاسخ آنلاین بستری را فراهم می‌کنند که کاربران بتوانند دانش خود را در حوزه‌های مختلف به اشتراک بگذارند و به دنبال پاسخ‌های مرتبط با سوالات خود باشند. این پلتفرم‌ها به دلیل توانایی‌شان در ارائه پاسخ‌های تخصصی و متناسب با نیاز کاربران، به یکی از منابع اصلی اطلاعات تبدیل شده‌اند. با این حال،

با رشد سریع جوامع پرسش و پاسخ و افزایش چشمگیر حجم اطلاعات، نیاز به سیستم‌های هوشمند برای تشخیص خودکار بهترین پاسخ‌ها بیش از پیش احساس می‌شود. در عصر دیجیتال، تعاملات انسانی به طور گسترده به فضای آنلاین منتقل شده است و یکی از مهم‌ترین این فضاها، جوامع پرسش و پاسخ آنلاین هستند. این پلتفرم‌ها، با ایجاد بستری برای طرح پرسش‌ها و دریافت پاسخ‌های تخصصی، به ابزاری قدرتمند برای حل مسائل و ارتقای دانش فردی و جمعی تبدیل شده‌اند.

با این حال، این جوامع با چالش‌هایی نظیر حجم عظیم اطلاعات، تنوع موضوعات و کیفیت متغیر پاسخ‌ها مواجه هستند که می‌توانند کارایی آن‌ها را محدود کنند. در چنین شرایطی، بهره‌گیری از روش‌های پیشرفته هوش مصنوعی و یادگیری ماشین به عنوان راهکاری مؤثر مطرح می‌شود.

این پژوهش با هدف بهبود فرآیند شناسایی بهترین پاسخ‌ها و ارتقای کارایی این پلتفرم‌ها انجام شده است. در این راستا، مدل‌های قدرتمند زبانی نظیر *XLNet* و

غیرممتنی پاسخ‌ها، به کاربران در یافتن پاسخ‌های بهینه کمک کرد. این کار نه تنها باعث افزایش کارایی و دقت در شناسایی پاسخ‌های مناسب می‌شود، بلکه تجربه کاربری را نیز بهبود می‌بخشد. از سوی دیگر، با توجه به نقش مهم اطلاعات در تصمیم‌گیری‌های روزمره و حرفه‌ای، ارائه روشی کارآمد برای بهبود دسترسی به پاسخ‌های باکیفیت، ضرورتی غیرقابل انکار است.

علاوه بر این، توسعه مدلی که بتواند ویژگی‌های متنی و غیرممتنی را به صورت هم‌زمان تحلیل کند، می‌تواند به پیشرفت‌های قابل توجهی در حوزه‌های مرتبط با پردازش زبان طبیعی و یادگیری ماشین منجر شود. این پژوهش با تمرکز بر داده‌های واقعی و روش‌های نوین، می‌تواند به عنوان پایه‌ای برای توسعه سیستم‌های توصیه‌گر در جوامع پرسش و پاسخ عمل کرده و نقش مهمی در بهبود تعاملات کاربران ایفا کند.

## ۲-۲- چالش‌های موجود

با وجود رشد روزافزون پلتفرم‌های پرسش و پاسخ آنلاین و اهمیت عملیاتی این سیستم‌ها، پیش‌بینی خودکار بهترین پاسخ همچنان با چالش‌های متعددی روبرو است. این چالش‌ها از ویژگی‌های ساختاری و مفهومی این جوامع سرچشمه می‌گیرند و نیازمند رویکردهایی پیشرفته و ترکیبی برای حل مسائل موجود هستند.

یکی از چالش‌های اساسی، نبود پاسخ‌های واضح و قطعی برای بسیاری از سوالات است. در بسیاری از موارد، پاسخ‌دهندگان از دیدگاه‌ها و رویکردهای مختلف به سوال پاسخ می‌دهند، که این موضوع انتخاب بهترین پاسخ را دشوار می‌کند. علاوه بر این، تعداد زیاد شرکت‌کنندگان در بحث و تعاملات پیچیده میان کاربران باعث می‌شود که ارزیابی و تحلیل کیفی پاسخ‌ها پیچیدگی بیشتری پیدا کند.

بسیاری از پژوهش‌های قبلی این مسائل را به طور کامل مدنظر قرار نداده‌اند. برای مثال، در سیستم‌های پرسش و

چالش‌های متعددی در این زمینه وجود دارد، از جمله حجم زیاد اطلاعات، کیفیت متغیر پاسخ‌ها و دشواری در شناسایی پاسخ‌های بهینه.

مسئله اصلی این است که کاربران غالباً نمی‌توانند به راحتی بهترین پاسخ ممکن را از میان تعداد زیادی از پاسخ‌های ارائه‌شده پیدا کنند. این موضوع منجر به صرف زمان زیاد و کاهش کارایی در استفاده از این جوامع می‌شود. علاوه بر این، معیارهای ارزیابی پاسخ‌ها به دلیل تأثیرگذاری عوامل مختلف متنی و غیرممتنی همچنان موضوعی چالش‌برانگیز است.

در این پژوهش، هدف ارائه راهکاری مبتنی بر هوش مصنوعی است که بتواند با تحلیل ویژگی‌های متنی (مانند کیفیت محتوا) و غیرممتنی (مانند تعداد رأی‌های مثبت) بهترین پاسخ را به صورت خودکار شناسایی کند. تمرکز این مطالعه بر شناسایی عواملی است که یک پاسخ را به عنوان بهترین پاسخ برای کاربران تبدیل می‌کند، به طوری که نیاز اطلاعاتی آن‌ها را به بهترین شکل ممکن برآورده سازد.

## ۲-۱- ضرورت انجام پژوهش

در جوامع پرسش و پاسخ، تعاملات کاربران و تنوع محتوای تولید شده پیچیدگی زیادی دارد و الگوریتم‌های معمول اغلب در مدل‌سازی این پیچیدگی‌ها و پاسخ‌گویی به نیازهای اطلاعاتی کاربران ناکارآمد هستند.

با افزایش استفاده از پلتفرم‌های پرسش و پاسخ آنلاین، این فضاها به یکی از منابع اصلی برای جستجوی اطلاعات در موضوعات متنوع تبدیل شده‌اند. با این حال، حجم بالای داده‌ها و تنوع کیفیت پاسخ‌ها، چالشی جدی برای کاربران به شمار می‌رود. شناسایی و انتخاب بهترین پاسخ در میان حجم انبوهی از اطلاعات موجود، نه تنها زمان‌بر است بلکه گاهی منجر به تصمیم‌گیری‌های نادرست می‌شود.

اهمیت پژوهش حاضر در این است که بتوان با بهره‌گیری از روش‌های هوش مصنوعی و تحلیل ویژگی‌های متنی و

- **سوگیری کاربران:** نظرات ذهنی کاربران و الگوهای رفتاری آن‌ها می‌تواند به نتایجی نامعتبر منجر شود.
- **پاسخ‌های چندمعنایی و متناقض:** وجود پاسخ‌هایی با تفسیرهای مختلف و گاه متناقض، شناسایی پاسخ بهینه را دشوار می‌سازد.
- **عدم تفسیرپذیری نتایج:** مدل‌های پیشرفته یادگیری ماشین گاهی فاقد شفافیت کافی برای توضیح انتخاب‌های خود هستند.

این چالش‌ها نشان می‌دهند که توسعه سیستم‌های پیشرفته‌تر برای پیش‌بینی بهترین پاسخ، نه تنها نیازمند استفاده از تکنیک‌های یادگیری ماشین و پردازش زبان طبیعی است، بلکه به ترکیب روش‌های چندگانه و بهره‌گیری از داده‌های متنی و غیرمتنی نیز نیاز دارد.

برای رفع این محدودیت‌ها، پژوهش حاضر تلاش دارد با ارائه رویکردی ترکیبی که ویژگی‌های زبانی، غیرمتنی و رابطه‌ای را در نظر می‌گیرد، دقت و کارایی انتخاب خودکار بهترین پاسخ را افزایش دهد. هدف نهایی این است که با توسعه مدل‌های هوشمند و قابل تفسیر، پلتفرم‌های پرسش و پاسخ آینده قادر باشند حتی در شرایط پیچیده و با وجود عدم قطعیت، پاسخ‌های بهینه را به کاربران ارائه دهند.

## ۲-۳- مفروضات و سوال‌های پژوهش

این پژوهش بر اساس مفروضاتی طراحی شده است که نقش اساسی در شکل‌گیری روش پیشنهادی ایفا می‌کنند. نخست آنکه، فرض بر این است که استفاده از رویکرد آموزش نظارت‌شده می‌تواند با کاهش پیچیدگی مدل، دقت و کارایی آن را در شناسایی بهترین پاسخ‌ها افزایش دهد. همچنین، انتخاب تابع زیان آنتروپی متقاطع باینری به‌عنوان معیار بهینه‌سازی، نقش مهمی در بهبود فرایند یادگیری و دقت در تشخیص پاسخ‌های بهینه ایفا می‌کند. داده‌های مورد استفاده در این پژوهش از مجموعه پرسش و پاسخ‌های مرتبط با حوزه اندروید در

پاسخ حوزه باز، پاسخ‌ها اغلب در مجموعه داده‌های عمومی یا منابعی مانند ویکی‌پدیا قرار دارند که شامل پاسخ‌های غیرمتمرکز و نامحدود است. اما در سیستم‌های پرسش و پاسخ اجتماعی، ساختار ساده‌تر است؛ سوال توسط یک کاربر مطرح شده و اعضای جامعه به آن پاسخ می‌دهند. با این حال، نبود معیارهای مشخص و جامع برای ارزیابی کیفیت پاسخ‌ها همچنان یک معضل کلیدی است.

یکی دیگر از چالش‌های برجسته، **فرضیات ساده‌سازی شده در مدل‌های کنونی** است. این مدل‌ها غالباً فرض می‌کنند که کاربران سوال‌کننده دانش و دیدگاه مشابهی با پاسخ‌دهندگان دارند، در حالی که در عمل، تفاوت‌های قابل توجهی در سطح دانش، باورها و نیازهای اطلاعاتی کاربران وجود دارد. مدل‌های فعلی به دلیل نادیده گرفتن این تنوع، قادر به در نظر گرفتن تفاوت‌های مهم در فرایند انتخاب بهترین پاسخ نیستند.

چالش‌های کلیدی در پیش‌بینی بهترین پاسخ به شرح زیر هستند:

- **کمبود داده‌های برجسته‌خورده باکیفیت:** بسیاری از مجموعه داده‌های موجود فاقد اطلاعات کافی برای آموزش مدل‌های پیشرفته یادگیری ماشین هستند.
- **کیفیت پایین داده‌ها:** وجود داده‌های نویزی و پاسخ‌های بی‌کیفیت می‌تواند دقت مدل‌ها را به شدت کاهش دهد.
- **تعاملات پیچیده میان کاربران:** روابط متقابل کاربران، مانند رای‌دهی و بازخورد، نیازمند مدل‌سازی دقیق‌تری هستند.
- **تغییرات زمانی:** کیفیت و اهمیت پاسخ‌ها ممکن است در گذر زمان تغییر کند و مدل‌ها باید این تغییرات را منعکس کنند.
- **تنوع زبانی و فرهنگی:** تفاوت در شیوه بیان و درک مفاهیم میان کاربران با پیشینه‌های مختلف می‌تواند فرایند پیش‌بینی را پیچیده‌تر کند.

مجموعه‌ای از رویکردهای نوآورانه و تکنیک‌های پیشرفته را به کار گرفته است.

#### ۲-۴-۱- استفاده از مدل‌های پیشرفته زبانی

این پژوهش از مدل‌های زبانی قدرتمند **XLNet** و **RoBERTa** برای استخراج خودکار ویژگی‌های متنی استفاده کرده است. این مدل‌ها با توانایی بی‌نظیر در درک معنایی عمیق و تحلیل پیچیدگی‌های متنی، نقش کلیدی در بهبود عملکرد مدل نهایی ایفا کرده‌اند.

#### ۲-۴-۲- بهره‌گیری از یادگیری جمعی وزن‌دار پشته‌ای

مدل پیشنهادی با استفاده از تکنیک استکینگ وزن‌دار به شکلی کاملاً هوشمندانه، خروجی‌های مدل‌های پایه مانند **XLNet**، **RoBERTa** و جنگل‌های تصادفی را ترکیب کرده است. این روش از نقاط قوت هر مدل بهره‌برداری کرده و به شکل مؤثری ضعف‌های آن‌ها را پوشش داده است.

#### ۲-۴-۳- استفاده از خروجی درصد وقوع برای پیش‌بینی نهایی

به جای استفاده از پیش‌بینی‌های قطعی صفر و یک، از خروجی درصد وقوع کلاس‌ها برای تصمیم‌گیری نهایی استفاده شده است. این رویکرد امکان تحلیل دقیق‌تر و انعطاف‌پذیری بیشتری در تعیین بهترین پاسخ‌ها فراهم کرده است.

#### ۲-۴-۴- ترکیب هوشمندانه ویژگی‌های متنی و غیرمتنی

به جای استفاده از لایه‌های کاملاً متصل سنتی، ویژگی‌های متنی استخراج‌شده از مدل‌های زبانی و ۲۱ ویژگی غیرمتنی دست‌ساز به شکل هوشمندانه‌ای ترکیب شده‌اند. این ترکیب کارآمد باعث بهبود دقت مدل و بهره‌برداری بهتر از اطلاعات ساختاریافته و معنایی مجموعه داده شده است.

#### ۲-۴-۵- مقیاس وسیع داده‌ها

این پژوهش از ۳۵,۵۶۸ جفت پرسش و پاسخ استخراج‌شده

پلتفرم **StackOverflow**، یکی از زیرمجموعه‌های شبکه **StackExchange**، جمع‌آوری شده است. این داده‌ها شامل ویژگی‌های متنی و غیرمتنی هستند که به صورت جداگانه استخراج و سپس ترکیب می‌شوند تا تأثیر استفاده از ویژگی‌های چندگانه در بهبود عملکرد مدل بررسی شود.

در راستای تحقق اهداف پژوهش، سؤالاتی مطرح می‌شود که جهت‌گیری تحقیق را مشخص می‌کنند. یکی از این پرسش‌ها این است که آیا رویکردهای کلاس‌بندی مبتنی بر برجسبدهی می‌توانند استراتژی مناسبی برای آموزش مدل در شناسایی بهترین پاسخ‌ها باشند. همچنین این سؤال بررسی می‌شود که آیا استفاده از یادگیری عمیق به صورت انتها به انتها می‌تواند استخراج ویژگی‌های متنی را به شکلی دقیق‌تر و اثربخش‌تر انجام دهد. از سوی دیگر، تأثیر ترکیب ویژگی‌های متنی و غیرمتنی بر دقت و کارایی مدل مورد توجه قرار می‌گیرد. در نهایت، این پژوهش بررسی می‌کند که آیا استفاده از مجموعه داده **StackOverflow**، به عنوان منبعی تخصصی در حوزه اندروید، می‌تواند در شناسایی دقیق‌تر پرسش‌ها و پاسخ‌های مرتبط مؤثر باشد.

#### ۲-۴- نوآوری پژوهش

پیش‌بینی بهترین پاسخ در جوامع پرسش و پاسخ به عنوان یکی از مسائل مهم در حوزه هوش مصنوعی و یادگیری ماشین شناخته می‌شود. در این پژوهش، ما با استفاده از یک مدل یادگیری جمعی پشته‌ای وزن‌دار، تحولی جدید در دقت و کارایی این سیستم‌ها ایجاد کرده‌ایم. مدل پیشنهادی ما با ترکیب چندین روش مدرن و نوآورانه، توانسته است دقت پیش‌بینی را به میزان ۹۹.۰۳٪ افزایش دهد، که در مقایسه با پژوهش‌های پیشین، یک جهش چشمگیر محسوب می‌شود.

پژوهش حاضر با هدف ارائه راهکاری جامع و دقیق برای پیش‌بینی بهترین پاسخ‌ها در جوامع پرسش و پاسخ،

برای جلوگیری از بیش‌برازش<sup>۲</sup> و افزایش قابلیت تعمیم‌دهی مدل، از تکنیک اعتبارسنجی متقاطع<sup>۳</sup> استفاده شده است. این روش به‌ویژه در شرایطی که داده‌های برچسب‌خورده محدود هستند، نقش مؤثری در مدیریت این محدودیت ایفا کرده و عملکرد مدل را بهبود بخشیده است. همچنین، مدل پیشنهادی با بهره‌گیری از ترکیب ویژگی‌های متنی (مانند معنای عمیق متن) و غیرمتنی (مانند تاریخچه تعاملات کاربران) توانسته ارزیابی جامعی از کیفیت پاسخ‌ها ارائه دهد.

یکی دیگر از ویژگی‌های منحصربه‌فرد این مدل، استفاده از فرایادگیرنده درخت تصمیم است که با دریافت خروجی مدل‌های پایه و بهره‌گیری از درصد وقوع کلاس‌ها به جای پیش‌بینی کلاس دقیق، ترکیب هوشمندانه‌ای از پیش‌بینی‌ها ایجاد کرده است. این روش علاوه بر دقت بالا، انعطاف‌پذیری مدل را نیز افزایش داده است.

در نهایت، مدل پیشنهادی توانایی قابل‌توجهی در مدیریت داده‌های کم دارد. استفاده از اعتبارسنجی متقاطع و ترکیب هوشمندانه مدل‌های پایه به کمک فرایادگیرنده، این امکان را فراهم کرده است که مدل در شرایط واقعی و با داده‌های محدود نیز عملکرد مناسبی داشته باشد. این ویژگی‌ها مدل پیشنهادی را به یک چارچوب قدرتمند و قابل اعتماد برای پیش‌بینی دقیق و انعطاف‌پذیر تبدیل کرده است.

### ۳- پژوهش‌های پیشین

پژوهشی که توسط Mamykina و همکاران [۱] (۲۰۱۱) انجام شده است، عوامل موفقیت سایت

android.stackexchange.com استفاده

کرده است. این حجم بالای داده‌ها، امکان آموزش و ارزیابی دقیق مدل را فراهم کرده و تعمیم‌پذیری نتایج را افزایش داده است.

### ۶-۴-۲- دستیابی به دقت بی‌نظیر

مدل نهایی با دقت ۹۹.۰۳٪، توانسته است عملکردی فوق‌العاده در پیش‌بینی بهترین پاسخ‌ها ارائه دهد. این دقت نشان‌دهنده برتری مدل پیشنهادی نسبت به رویکردهای پیشین و توانایی بالای آن در حل مسئله است.

این نوآوری‌ها در کنار یکدیگر، پژوهش حاضر را به یک راهکار جامع و پیشرفته برای بهبود کیفیت تعاملات در جوامع پرسش و پاسخ تبدیل کرده‌اند.

### ۵-۲- مزایای کلیدی مدل پیشنهادی

مدل پیشنهادی با استفاده از ترکیبی از مدل‌های پایه قوی مانند XLNet، RoBERTa، جنگل‌های تصادفی و رأی‌گیری اکثریت طراحی شده است که هر کدام به‌طور خاص در تحلیل ویژگی‌های متنی و غیرمتنی تخصص دارند. این تنوع مدل‌ها نه تنها به یادگیری بهتر کمک کرده بلکه باعث افزایش انعطاف‌پذیری و بهبود عملکرد مدل نیز شده است. یکی از نوآوری‌های کلیدی این پژوهش استفاده از تکنیک وزن‌دار پشته‌ای<sup>۱</sup> است، جایی که وزن هر مدل پایه بر اساس دقت و عملکرد آن به‌دقت تنظیم شده است. این تکنیک با کاهش خطاهای مدل‌های پایه توانسته دقت پیش‌بینی نهایی را به طور چشمگیری افزایش دهد.

<sup>۳</sup> Cross-Validation

<sup>۱</sup> Weighted Stacking

<sup>۲</sup> Overfitting

در پژوهش Rajpurkar و همکاران (۲۰۱۶) [۳]، مجموعه داده Stanford Question Answering Dataset (SQuAD) معرفی شده است که شامل بیش از ۱۰۰,۰۰۰ پرسش و پاسخ از مقالات ویکی پدیا است. این مجموعه به عنوان چالشی برای مدل های درک مطلب طراحی شده است و تنوع بالایی از پرسش ها و انواع پاسخ ها را شامل می شود. پژوهش نشان داد که عملکرد انسانی (۸۶.۸٪) به مراتب بالاتر از بهترین مدل (۵۱٪) است، به ویژه با افزایش پیچیدگی پاسخ ها و تفاوت های نحوی میان سوال و جمله حاوی پاسخ. اگرچه مدل رگرسیون لجستیک برای برخی از انواع پاسخ ها عملکرد قابل قبولی داشت، اما در تشخیص دقیق محدوده پاسخ ها ضعف نشان داد. محدودیت های این پژوهش شامل تمرکز صرف بر مقالات ویکی پدیا و عدم توانایی مدل ها در مدیریت کامل تفاوت های نحوی است. بهبود مدل های پردازش زبان و گسترش مجموعه داده به موضوعات متنوع تر از پیشنهاد های آتی پژوهش است.

در سال های اخیر، توسعه سیستم های پرسش و پاسخ (QA) مبتنی بر پایگاه های دانش مورد توجه بسیاری از پژوهش ها قرار گرفته است. برای مثال، Cui و همکاران (۲۰۱۹) [۴] سیستم KBQA را معرفی کردند که از الگوهای یادگیری برای تبدیل سوالات طبیعی به پرسش های ساختاریافته استفاده می کند. این سیستم توانست با گسترش روابط در پایگاه دانش، پوشش پایگاه را تا ۵۷ برابر افزایش دهد و دقت بالایی در پاسخ به سوالات پیچیده به دست آورد. KBQA در مقایسه با روش های پیشرفته دیگر، در معیار QALD، دقت و بازخوانی بیشتری داشت و در شناسایی موجودیت ها نیز دقت ۷۲٪ را ثبت کرد. با این حال، عملکرد آن در پاسخ به سوالات غیردوتایی محدود بود که باعث کاهش بازخوانی و دقت کلی در برخی موارد شد.

در سال های اخیر، پژوهش های متعددی به بهبود عملکرد سامانه های پاسخ دهی به سوالات در حوزه های باز

پرسش و پاسخ Stack Overflow را بررسی می کند. این مطالعه نشان داد که موفقیت این پلتفرم به طراحی هوشمندانه، تعامل فعال بنیان گذاران با جامعه کاربری، و یک سیستم تشویقی مبتنی بر رفتارهای مطلوب بستگی دارد. نتایج تحقیق نشان می دهد که سایت دارای نرخ پاسخ دهی بیش از ۹۰ درصد و زمان میانگین ۱۱ دقیقه برای دریافت اولین پاسخ است. پژوهش همچنین به اکولوژی منحصر به فرد رفتار کاربران در این پلتفرم اشاره دارد که شامل گروه های مختلفی مانند فعالان جامعه، کاربران پرکار مقطعی، کاربران کم فعال و بازدیدکنندگان است. اگرچه این طراحی منجر به ایجاد یک منبع پایدار برای توسعه دهندگان نرم افزار شده است، محدودیت هایی نظیر مشکل «سریع ترین پاسخ دهنده» و تفاوت های میان کاربران تازه کار و خبره نیز شناسایی شده است. این پژوهش بر اهمیت طراحی پاسخ گو، به روزرسانی مداوم، و تعامل مستمر با کاربران برای موفقیت پلتفرم های اجتماعی تأکید می کند.

پژوهشی توسط Gkotsis و همکاران (۲۰۱۵) [۲] سیستم نوآورانه ای به نام ACQUA را برای پیش بینی بهترین پاسخ در وبسایت های پرسش و پاسخ معرفی می کند. این سیستم از ویژگی های متنی ساده برای ارزیابی پاسخ ها استفاده کرده و توانسته است عملکردی بهتر از روش های مبتنی بر زبان شناسی و معادل با روش های مبتنی بر امتیازدهی کاربران ارائه دهد. نتایج نشان می دهد که دقت میانگین سیستم ۸۴٪ و نرخ بازیابی آن ۷۰٪ است. این مطالعه همچنین نشان می دهد که گسسته سازی ویژگی های زبانی باعث بهبود چشمگیر عملکرد می شود، در حالی که ویژگی هایی مانند شهرت کاربران تأثیر قابل توجهی بر پیش بینی بهترین پاسخ ندارند. این تحقیق رویکردی موثر برای شناسایی پاسخ های مناسب در جوامع آنلاین ارائه کرده و بر کاربرد ویژگی های ساده متنی برای ارزیابی محتوای تولید شده توسط کاربران تأکید دارد.

بی کیفیت و نیاز به بهبود دقت مدل پیش‌بینی است. این پژوهش کاربردهایی در بهبود کیفیت محتوای پلتفرم‌های مشابه از طریق بازخورد فوری به کاربران و حمایت از مدیران دارد.

در پژوهش Jeon و همکاران (۲۰۰۶) [۷]، چارچوبی برای پیش‌بینی کیفیت پاسخ‌ها با استفاده از ویژگی‌های غیرمتمنی توسعه داده شد. این مطالعه از روش‌هایی مانند "تخمین چگالی هسته‌ای" و "مدل آنتروپی بیشینه" برای پردازش داده‌ها استفاده کرده و نشان داده است که ترکیب این ویژگی‌ها در مدل‌های بازیابی اطلاعات می‌تواند عملکرد بازیابی پرسش‌های متداول را به طور قابل توجهی بهبود بخشد. آزمایش‌ها نشان داد که پیش‌بینی‌کننده توسعه‌یافته توانست پاسخ‌های باکیفیت را با دقت میانگین ۰.۹۲۲۷ و پاسخ‌های بی کیفیت را با دقت ۰.۶۵۵۸ از یکدیگر تفکیک کند. علاوه بر این، با تلفیق این معیار کیفیت در مدل‌های بازیابی، عملکرد رتبه‌بندی به طور معناداری بهبود یافت (با سطح اطمینان ۰.۹۹٪). با این حال، محدودیت‌هایی مانند حجم کوچک داده‌های آموزشی (۸۹۴ نمونه آموزشی) و نیاز به توسعه معیارهای استاندارد برای سنجش کیفیت مستندات در پژوهش مشخص شد.

در پژوهش Shah و Pomerantz (۲۰۱۰) [۸]، کیفیت پاسخ‌ها در پلتفرم‌های پرسش و پاسخ اجتماعی مانند Yahoo! Answers ارزیابی و پیش‌بینی شد. این مطالعه دو رویکرد را برای سنجش کیفیت پاسخ‌ها مقایسه کرد: ارزیابی انسانی بر اساس ۱۳ معیار مشخص و استخراج خودکار ویژگی‌ها از داده‌ها. نتایج نشان داد که ویژگی‌های استخراج‌شده خودکار مانند نمایه کاربران و رتبه‌بندی متقابل پاسخ‌ها نسبت به معیارهای انسانی کارایی بهتری در پیش‌بینی پاسخ‌های باکیفیت دارند و دقت پیش‌بینی را از ۸۰.۳۳٪ به ۸۴.۱۷٪ افزایش دادند. با این حال، مدل لجستیک رگرسیون با استفاده از معیارهای انسانی نتوانست به‌طور موثر متغیرهای داده را توضیح دهد (با  $\text{pseudo } R^2$  برابر با ۰.۰۹۰۲). پژوهش همچنین بر

پرداخته‌اند. Zhou و همکاران (۲۰۲۰) [۵] روش پاسخ‌دهی به سوالات حوزه باز مبتنی بر دانش (KAQA) را معرفی کردند که با استفاده از سه‌گانه‌های دانشی استخراج‌شده از منابع خارجی، روابط میان سوالات و اسناد و همچنین روابط میان اسناد را مدل‌سازی می‌کند. این چارچوب شامل سه مؤلفه اصلی است: بازیاب، خواننده و رتبه‌بند مجدد. آزمایش‌ها نشان داد که KAQA دقت بازیابی و رتبه‌بندی پاسخ‌ها را در مجموعه داده‌هایی نظیر SQuAD-open، Quasar-T و TriviaQA بهبود داده و عملکرد کلی سیستم‌های پاسخ‌دهی را افزایش می‌دهد. با این حال، این روش محدود به استفاده از مجموعه داده‌های مشخص بوده و از مدل‌های پیچیده‌تر مانند شبکه‌های گراف کانولوشنی برای گسترش دانش رابطه‌ای بهره نمی‌برد. این روش از بازیاب و رتبه‌بند مجدد مبتنی بر گراف‌های سوال-سند و سند-سند استفاده می‌کند که باعث بهبود دقت بازیابی و رتبه‌بندی پاسخ‌ها شده است. هرچند میزان دقت دقیق (Exact Match یا F۱) گزارش نشده است، اما تاکید پژوهشگران بر بهبود عملکرد کلی نسبت به سایر روش‌های موجود نشان‌دهنده تاثیر قابل توجه آن بر دقت سیستم‌های پاسخ‌دهی حوزه باز است.

در پژوهش Correa و Sureka (۲۰۱۴) [۶]، رفتار و ویژگی‌های سوالات حذف‌شده در Stack Overflow بررسی شده و چارچوبی پیش‌بینی‌کننده برای شناسایی سوالات باکیفیت پایین هنگام ایجاد پیشنهاد شده است. این مطالعه با تحلیل داده‌های پنج‌ساله و استفاده از ۴۷ ویژگی از جمله نمایه کاربر، محتوای سوال و سبک نحوی، مدلی با دقت ۶۶٪ در پیش‌بینی سوالات حذف‌شده ارائه کرده است. نتایج نشان داد که حدود ۸٪ از سوالات حذف می‌شوند، که اغلب توسط مدیران و پس از زمان قابل توجهی از دریافت اولین رأی حذف صورت می‌گیرد. سوالات حذف‌شده عموماً کیفیت بسیار پایینی دارند و در پایین ساختار هرمی کیفیت سوالات قرار می‌گیرند. محدودیت‌ها شامل عدم بررسی انواع دیگر محتوای

BERTScore که بازنمایی‌های متنی کانتکت‌محور را به کار می‌گیرد، بهبود جزئی نسبت به معیارهای سنتی نشان داد اما همچنان نیاز به بهینه‌سازی بیشتر دارد. این پژوهش به محدودیت معیارهای موجود در درک پاسخ‌های آزاد و پیچیده اشاره کرد و بر ضرورت توسعه ابزارهای دقیق‌تر برای ارزیابی سیستم‌های پاسخ به پرسش تأکید داشت.

Hu و همکاران (۲۰۱۷) [۱۱] رویکردی مبتنی بر یادگیری عمیق برای پیش‌بینی کیفیت پاسخ‌های خدمات پرسش و پاسخ آنلاین در حوزه سلامت پیشنهاد کردند که ویژگی‌های متنی و غیرمتنی را با استفاده از یک چارچوب شبکه باور عمیق چندرسانه‌ای<sup>۴</sup> ترکیب می‌کند. نتایج نشان داد که ویژگی‌های غیرمتنی، شامل ویژگی‌های زبانی سطحی و ویژگی‌های اجتماعی، نقش مؤثری در تمایز پاسخ‌های با کیفیت از پاسخ‌های کم‌کیفیت دارند. ویژگی‌های زبانی سطحی عملکرد بهتری نسبت به ویژگی‌های اجتماعی داشتند. این چارچوب نسبت به روش‌های پایه عملکرد بهتری نشان داد و دقت، فراخوانی و امتیاز  $F_1$  بالاتری داشت. با این حال، مطالعه محدود به مجموعه داده خاصی از پلتفرم Haodf Online بود و ممکن است نتایج آن به سایر دامنه‌ها تعمیم‌پذیر نباشد.

Agichtein و همکاران (۲۰۰۸) [۱۲] چارچوبی برای ارزیابی کیفیت محتوای تولیدشده توسط کاربران در رسانه‌های اجتماعی، به‌ویژه در پورتال‌های پرسش و پاسخ جامعه‌ای مانند Yahoo! Answers، معرفی کردند. این رویکرد با ترکیب ویژگی‌های مبتنی بر محتوا، روابط کاربران و استفاده از داده‌ها (مانند بازدید صفحات)، به شناسایی محتوای با کیفیت بالا پرداخت. نتایج آزمایش‌ها

اهمیت اطلاعات اجتماعی و زمینه‌ای در ارزیابی کیفیت محتوا تأکید کرد. محدودیت‌های این مطالعه شامل کوچک بودن حجم داده‌ها و سوگیری احتمالی در ارزیابی انسانی است.

در سال‌های اخیر، بسیاری از پژوهش‌ها به بهبود شناسایی خودکار پست‌های کم‌کیفیت در پلتفرم‌های پرسش و پاسخ پرداخته‌اند. برای مثال، Ponzanelli و همکاران (۲۰۱۴) [۹] با استفاده از ویژگی‌های متنی و اجتماعی، مدلی مبتنی بر الگوریتم‌های ژنتیک ارائه دادند که توانست صف بررسی پست‌های کم‌کیفیت در Stack Overflow را کاهش داده و پست‌های باکیفیت اشتباه‌شناسایی‌شده را با حداقل خطا حذف کند. این رویکرد با دستیابی به دقت نرم ۶۸.۹۵٪ و کاهش ۹ درصدی صف بررسی، زمان بررسی توسط مدیران را بهبود بخشید. با این حال، محدودیت داده‌های عمومی موجود از Stack Overflow ممکن است بر دقت نتایج تأثیرگذار باشد. پیشنهاد شده است که این روش برای سیستم‌های توصیه‌گر ساده یا بهبود خودکار صف بررسی توسعه یابد.

همچنین، ارزیابی سیستم‌های پاسخ به پرسش توجه زیادی را به خود جلب کرده است. برای مثال، Chen و همکاران (۲۰۱۹) [۱۰] عملکرد معیارهای خودکار نظیر METEOR، ROUGE، BLEU و  $F_1$  را بر روی سه مجموعه داده‌ی NarrativeQA، ROPES و SemEval بررسی کردند. نتایج نشان داد که معیار METEOR در مجموعه داده‌های NarrativeQA و SemEval بیشترین تطابق را با ارزیابی‌های انسانی داشت، در حالی که  $F_1$  در مجموعه داده ROPES عملکرد ضعیفی ارائه داد. همچنین، استفاده از معیار

<sup>۴</sup> DBN



پیش‌بینی‌کننده‌های کیفیت پاسخ استفاده شوند. مدل توسعه‌یافته دارای ضریب R برابر ۰.۸۸۶ بود که رابطه قوی بین مدل و کیفیت پاسخ را نشان می‌دهد. محدودیت‌های این مطالعه شامل استفاده از یک مجموعه داده خاص از Yahoo! Answers و محدودیت به یک دسته‌بندی خاص (علوم کامپیوتر) و نمونه‌گیری کوچک (۳۰۰ جفت سؤال و پاسخ) است. به علاوه، این مطالعه تنها بر ویژگی‌های متنی تمرکز داشت و ویژگی‌های غیرمتنی تأثیر کمی در کیفیت پاسخ داشتند. نتایج این مطالعه می‌توانند به بهبود مازول‌های استخراج پاسخ در سیستم‌های پرسش و پاسخ و رتبه‌بندی و بازیابی پاسخ‌ها کمک کنند. تحقیقات آینده می‌توانند این یافته‌ها را برای سیستم‌های پرسش و پاسخ خودکار و گسترده‌تر استفاده کنند و تحلیل‌های مشابه را برای داده‌های بیشتری از سایر دسته‌ها انجام دهند.

در مطالعه‌ای که توسط Suggu و همکاران (۲۰۱۶) [۱۴] انجام شد، یک معماری یادگیری عمیق به نام DFFN برای ارزیابی کیفیت پاسخ‌ها در سیستم‌های پرسش و پاسخ اجتماعی پیشنهاد شد. این معماری از ترکیب شبکه‌های عصبی کانولوشن و LSTM دوطرفه<sup>۷</sup> با توجه و ویژگی‌های دست‌ساز از منابع مختلف استفاده می‌کند. نتایج نشان داد که روش DFFN عملکردی برتر نسبت به روش‌های مرسوم دارد و در مجموعه داده‌های SemEval-۲۰۱۵ و SemEval-۲۰۱۶ نتایج بهتری را ارائه کرد. مدل پیشنهادی به‌ویژه DFFN-BLNA، بیشترین عملکرد را از خود نشان داد. این رویکرد با ترکیب ویژگی‌های استخراج‌شده از منابع خارجی مانند متن‌های ویکی‌پدیا و داده‌های کلیک، توانست ویژگی‌های تشابهی

نشان داد که این سیستم می‌تواند محتوای با کیفیت را با دقت مشابه ارزیابی‌های انسانی شناسایی کند. در ارزیابی عملکرد، معیارهای دقت<sup>۵</sup>، فراخوانی<sup>۶</sup> و مساحت زیر منحنی ROC (AUC) گزارش شد که نشان‌دهنده عملکرد بالا و دقت مناسب مدل بود. مهم‌ترین ویژگی‌هایی که در این مطالعه شناسایی شدند شامل طول پاسخ‌ها، دسته‌بندی موضوعی و آمار استفاده از محتوا (مانند تعداد بازدید صفحات) بودند. همچنین، این مطالعه نشان داد که ویژگی‌های مختلف مانند روابط کاربران و آمار استفاده می‌توانند به دقت بیشتری در شناسایی محتوای با کیفیت کمک کنند. با این حال، محدودیت‌هایی وجود داشت؛ این روش به طور خاص برای پورتال‌های پرسش و پاسخ طراحی شده و ممکن است به سایر حوزه‌ها تعمیم‌پذیر نباشد. همچنین، شامل چالش‌هایی مانند بیش‌برازش در صورتی که تنها از آمار استفاده بدون در نظر گرفتن بازخوردهای دیگر استفاده شود، بود.

Blooma و همکاران (۲۰۰۸) [۱۳] در این مطالعه چارچوبی پیش‌بینی‌کننده برای ارزیابی کیفیت پاسخ‌ها در سیستم‌های پرسش و پاسخ مبتنی بر جامعه ارائه دادند. این چارچوب از ویژگی‌های متنی و غیرمتنی برای پیش‌بینی بهترین پاسخ استفاده می‌کند. نتایج نشان داد که ویژگی‌های متنی مانند دقت، تکمیل و منطقی بودن پاسخ، مهم‌ترین پیش‌بینی‌کننده‌ها برای کیفیت پاسخ بودند. ویژگی‌های غیرمتنی مانند شهرت و اعتبار پاسخ‌دهنده یا پرسش‌کننده تأثیر چندانی نداشتند. مطالعه با استفاده از یک تحلیل رگرسیون برای مدلسازی رابطه میان ویژگی‌ها و کیفیت پاسخ‌ها انجام شد. نتایج نشان داد که دقت، تکمیل و ارتباط پاسخ‌ها می‌توانند به‌عنوان

<sup>۷</sup> BLSTM<sup>۵</sup> Precision<sup>۶</sup> Recall

Gkotsis و همکاران (۲۰۱۵) [۱۶] روشی نوین برای پیش‌بینی بهترین پاسخ در پلتفرم‌های پرسش و پاسخ مبتنی بر جامعه با استفاده از ویژگی‌های زبانی سطحی گسسته‌شده ارائه کردند. روش آن‌ها از روش‌های سنتی که به رتبه‌بندی و شهرت کاربران تکیه دارند، پیشی گرفت. این مطالعه نشان داد که استفاده از این ویژگی‌های گسسته‌شده، که شامل گروه‌بندی و مرتب‌سازی ویژگی‌های زبانی است، دقت بالاتری به همراه دارد و با F-Measure که یک معیار برای ارزیابی مدل‌های طبقه‌بندی است و ترکیبی از دقت و یادآوری است، برابر با ۰.۷۷ و AUC برابر با ۰.۸۸ عملکرد بهتری را از خود نشان می‌دهد. با وجود این که روش پیشنهادی عملکرد بهتری نسبت به روش‌های سنتی داشت، این مطالعه اشاره کرد که ویژگی‌های زبانی به تنهایی محدودیت‌هایی دارند، به ویژه زمانی که با استفاده از رتبه‌بندی کاربران مقایسه می‌شود. همچنین، نتایج ممکن است به دلیل تفاوت‌های زبانی در سایت‌های مختلف StackExchange متغیر باشد. پژوهش‌های آینده به بررسی قابلیت اعمال این روش در زمینه‌های دیگر و ارزیابی تأثیر کیفیت متنی در انتخاب پاسخ‌ها پیشنهاد شده است.

در سال‌های اخیر، پژوهش‌های مختلفی به ارزیابی کیفیت پاسخ‌ها در سایت‌های پرسش و پاسخ اجتماعی پرداخته‌اند. برای مثال، Zhu و همکاران (۲۰۰۹) [۱۷] مدلی چندبعدی برای ارزیابی کیفیت پاسخ‌ها معرفی کردند که شامل ۱۳ بعد مختلف مانند اطلاع‌رسانی، مودب بودن، تکمیل بودن، خوانایی، مرتبط بودن و صداقت است. این مدل از ترکیب تکنیک‌های کیفیت اطلاعات و پردازش زبان طبیعی استفاده می‌کند و قادر است فرآیند ارزیابی را به‌طور خودکار انجام دهد. نتایج اولیه نشان می‌دهد که این مدل می‌تواند دقت ۸۳.۹۸٪ را در پیش‌بینی پاسخ‌های خوب و بد بدست آورد. با این حال، برخی ابعاد مانند سطح جزئیات و تخصص نیاز به توسعه تکنیک‌های پردازش زبان طبیعی جدید دارند. این پژوهش‌ها به توسعه سیستم‌های ارزیابی کیفیت خودکار برای سایت‌های پرسش و پاسخ

را غنی‌تر کند. دقت مدل DFFN در مجموعه داده‌های SemEval-۲۰۱۶ و SemEval-۲۰۱۶ به طور خاص به صورت عددی ذکر نشده است، اما اشاره شده که این مدل عملکردی "state-of-the-art" داشت و توانست سیستم‌های پیشرفته‌ی دیگر را پشت سر بگذارد. بر اساس آزمایش‌های آماری انجام شده، نتایج به طور قابل توجهی بهتر از مدل‌های پایه بوده و در سطح اطمینان ۹۵٪ ( $\alpha = 0.05$ ) معنی‌دار بوده است. در نتیجه، مدل DFFN- BLNA بهترین عملکرد را در مقایسه با سایر مدل‌ها داشته است. در کل، نویسندگان ادعا کرده‌اند که این مدل به دقت بالاتر از سیستم‌های رقیب در همین مجموعه داده‌ها دست یافته است. این تحقیق همچنین محدودیتی در خصوص معایب یا کارهای آینده ارائه نکرده است. با وجود موفقیت‌های به‌دست‌آمده، یکی از چالش‌ها ممکن است نیاز به منابع داده‌ای غنی و زمان پردازش بالا باشد.

Hadfi و همکاران (۲۰۲۲) [۱۵] در مقاله خود به پیش‌بینی بهترین پاسخ در سایت‌های پرسش و پاسخ (CQA) با استفاده از اطلاعات کاربران پرداختند. در این مطالعه، ویژگی‌های مختلفی از جمله ویژگی‌های زبانی، اطلاعات کاربر، و ویژگی‌های رابطه‌ای برای پیش‌بینی بهترین پاسخ ترکیب شدند. نتایج نشان داد که روش پیشنهادی عملکرد بهتری نسبت به روش‌های قبلی داشته و ویژگی‌های زمان تاخیر بین سوال و پاسخ و امتیازدهی به پاسخ از توان پیش‌بینی بالاتری برخوردارند. در این تحقیق از مدل‌های مختلفی مانند SVM، Random Forest، MARS و LightGBM برای پیش‌بینی بهترین پاسخ استفاده شد که مدل LightGBM بهترین عملکرد را داشت. این روش به طور ویژه برای سایت‌های CQA مانند Stack Exchange مفید است، اما برای تعمیم یافته‌ها نیاز به ارزیابی‌های بیشتر با استفاده از داده‌های دیگر دارد. همچنین، محدودیت‌هایی مانند فرضیه جمعیت یکنواخت سوال‌کنندگان و حساسیت امتیازدهی به تغییرات پس از انتخاب بهترین پاسخ ذکر شده است.

عوامل موثر بر انتخاب بهترین پاسخ در سرویس‌های پرسش و پاسخ جامعه‌محور<sup>۱۱</sup> پرداختند. آن‌ها از رویکردی ترکیبی برای تحلیل ویژگی‌های اجتماعی، متنی و ارزیابی محتوا استفاده کردند و داده‌های خود را از Yahoo! Answers جمع‌آوری نمودند. نتایج تحلیل رگرسیون لجستیک نشان داد که ویژگی‌های ارزیابی محتوا، نظیر کامل بودن و دقت، تاثیر بیشتری در انتخاب بهترین پاسخ دارند. همچنین، ویژگی‌های متنی مانند طول پاسخ و نسبت طول پرسش و پاسخ نیز مؤثر بودند، اما نقش کمتری نسبت به ویژگی‌های محتوا داشتند. از سوی دیگر، ویژگی‌های اجتماعی مانند تأیید کاربران ارتباط ضعیف‌تری با انتخاب بهترین پاسخ نشان دادند. این پژوهش، محدود به یک دسته خاص (کامپیوتر و اینترنت) از داده‌های Yahoo! Answers بوده و ممکن است به سایر دسته‌ها یا سرویس‌های مشابه تعمیم‌پذیر نباشد. دقت خاصی به‌عنوان یک معیار کمی گزارش نشده است. نتایج این تحقیق بیشتر بر تحلیل اهمیت نسبی ویژگی‌های مختلف (اجتماعی، متنی و ارزیابی محتوا) در انتخاب بهترین پاسخ تمرکز داشت. اگرچه آن‌ها به خوبی نشان داده‌اند که ویژگی‌های ارزیابی محتوا بیشترین تأثیر را در این فرآیند دارند، اما دقت مدل یا چارچوب پیشنهادی برای پیش‌بینی بهترین پاسخ در داده‌ها بیان نشده است. علاوه بر این، حجم داده‌های مورد استفاده کوچک بوده و نتایج بر اساس ارزیابی‌های کارشناسی استخراج شده است. از این رو، مطالعه‌های آتی می‌توانند با استفاده از داده‌های گسترده‌تر و دسته‌بندی‌های متنوع‌تر، نقش سایر عوامل را بررسی کنند. همچنین، پیشنهاد شده است از روش‌هایی نظیر تحلیل لینک برای

اجتماعی کمک می‌کند، اما محدودیت‌هایی مانند استفاده از یک مجموعه داده خاص و نیاز به داده‌های بزرگتر برای بهبود دقت مدل وجود دارد.

Liu و همکاران (۲۰۱۰) [۱۸] مدلی برای پیش‌بینی بهترین پاسخ‌دهندگان در سایت‌های پرسش و پاسخ اجتماعی ارائه دادند. این مدل از ترکیب مدل زبان و تخصیص دیریکله نهفته<sup>۸</sup> برای مدل‌سازی پروفایل کاربران و پیش‌بینی پاسخ‌دهندگان استفاده می‌کند. با ترکیب فعالیت و اعتبار کاربران، مدل توانست دقت بالاتری در پیش‌بینی بهترین پاسخ‌دهندگان به دست آورد. نتایج تجربی با استفاده از معیار جدید "موفقیت در N" نشان داد که این روش می‌تواند به طور مؤثر سوالات جدید را به پاسخ‌دهندگان مناسب هدایت کند و کیفیت پاسخ‌ها را بهبود بخشد. با این حال، محدودیت‌هایی نظیر مقیاس کوچک داده‌ها و لزوم بهبود مدل‌سازی فعالیت و اعتبار کاربران وجود دارد. این روش کاربرد عملی در بهبود خدمات سایت‌های پرسش و پاسخ دارد، اما دقت پیش‌بینی می‌تواند به دلیل تعدد کاربران بالقوه کاهش یابد. دقت مدل به طور دقیق بیان نشده است، اما نتایج تجربی نشان داده‌اند که ترکیب مدل زبان و تخصیص دیریکله نهفته عملکرد بهتری نسبت به هر یک به تنهایی داشته است و افزودن اطلاعات مربوط به اعتبار کاربران و فعالیت کاربران نیز باعث بهبود دقت پیش‌بینی شده است. همچنین از معیاری به نام  $S@N^{10}$  برای ارزیابی عملکرد مدل استفاده شده است، که میزان موفقیت در توصیه بهترین پاسخ‌دهندگان را اندازه‌گیری می‌کند.

بلوم و همکاران (۲۰۱۰) [۱۹] در مطالعه‌ای به بررسی

<sup>۱۰</sup> Success-at-N<sup>۸</sup> LDA<sup>۱۱</sup> CQA<sup>۹</sup> Success-at-N

پیش‌بینی بهترین پاسخ در وبسایت‌های پرسش و پاسخ اجتماعی ارائه کردند که از ویژگی‌های زبانی سطحی گسسته‌شده استفاده می‌کند. این روش به دقت متوسط ۸۴٪ و بازخوانی ۷۰٪ دست یافت و عملکردی بهتر از روش‌های مبتنی بر زبان‌شناسی و مشابه با رویکردهای مبتنی بر امتیازدهی ارائه داد. در این پژوهش نشان داده شد که گسسته‌سازی ویژگی‌های زبانی میزان اطلاعات قابل استخراج را افزایش می‌دهد و استفاده از ویژگی‌های غیرزبانی تأثیری محدود بر عملکرد مدل دارد. این روش مستقل از داده‌های خاص هر جامعه بوده و بر روی ۲۱ وبسایت StackExchange آزمایش شده است. اگرچه مدل پیشنهادی برای شناسایی بهترین پاسخ مؤثر است، عدم استفاده از اطلاعات زمینه‌ای مانند امتیازات کاربران و پیچیدگی تغییرات زبانی در طول زمان محدودیت‌هایی ایجاد کرده است. پژوهش‌های آینده می‌توانند بر تحلیل تغییرات ویژگی‌های زبانی در طول زمان و استفاده از دیگر ویژگی‌ها برای بهبود عملکرد مدل تمرکز کنند.

Chen و همکاران (۲۰۱۲) [۲۲] روشی برای کالبراسیون رأی‌ها در سیستم‌های پرسش و پاسخ اجتماعی پیشنهاد کردند که با کاهش اثرات سوگیری کاربران در رأی‌گیری، برآورد کیفیت پاسخ و تخصص کاربران را بهبود می‌بخشد. این روش از یک مدل یادگیری نظارت‌شده استفاده می‌کند که وزن‌های اهمیت رأی‌ها را محاسبه کرده و کیفیت پاسخ‌ها را پیش‌بینی می‌کند. نتایج نشان می‌دهد که مدل کالیبره‌شده در تخمین کیفیت پاسخ و رتبه‌بندی تخصص کاربران نسبت به روش‌های غیرکالیبره بهبود قابل‌توجهی در معیارهایی مانند MAP و NDCG@k دارد. با این حال، کوتاه بودن فهرست رتبه‌بندی در ارزیابی محدودیتی در تمایز عملکرد

افزایش دقت و استفاده مجدد از پاسخ‌ها به‌عنوان منبع اطلاعاتی در سیستم‌های خودکار پرسش و پاسخ بهره گرفته شود.

Burel و همکاران (۲۰۱۲) [۲۰] مدلی برای شناسایی بهترین پاسخ‌ها در جوامع پرسش‌وپاسخ آنلاین طراحی کردند که از ویژگی‌های محتوا، کاربر و نخ‌های گفتگو بهره می‌برد. نتایج نشان داد ویژگی‌های نخ‌های گفتگو، مانند امتیازات و رتبه‌بندی‌ها، بیشترین تأثیر را در شناسایی بهترین پاسخ دارند، در حالی که ویژگی‌های کاربر، مانند نسبت شهرت موضوعی، نقش کمتری ایفا می‌کنند. ویژگی‌های نخ‌های گفتگو<sup>۱۲</sup> شامل اطلاعات ساختاری و تعاملی مربوط به یک موضوع در پلتفرم‌های پرسش‌وپاسخ است. این ویژگی‌ها شامل امتیازات و رأی‌ها (میزان رأی‌های مثبت و منفی)، تعداد پاسخ‌ها (کل پاسخ‌های ارسال‌شده در یک موضوع)، رتبه‌بندی پاسخ‌ها (مقایسه امتیاز یک پاسخ با دیگر پاسخ‌های همان موضوع)، زمان‌بندی‌ها (زمان ارسال پاسخ‌ها و ایجاد موضوع)، و ساختار نخ (ترتیب پاسخ‌ها و تعاملات کاربران) می‌شود. در مطالعه Burel و همکاران، این ویژگی‌ها تأثیر قابل‌توجهی در پیش‌بینی کیفیت پاسخ‌ها نشان داده و برتری آنها نسبت به ویژگی‌های مرتبط با محتوا و کاربران تأیید شده است. برخلاف پژوهش‌های قبلی، طول پاسخ با انتخاب بهترین پاسخ همبستگی معناداری نداشت. این مدل در سه مجموعه داده با دقت ۸۳٪ تا ۸۷٪ عملکرد موفقیت‌آمیزی داشت. با این حال، پژوهش محدود به سه جامعه آنلاین بوده و ممکن است به سایر پلتفرم‌ها قابل تعمیم نباشد.

Gkotsis و همکاران (۲۰۱۴) [۲۱] مدلی برای

طبقه‌بند متا، دقت، بازخوانی و F-Measure بهتری نسبت به روش‌های تکنما ارائه می‌دهد. نتایج نشان می‌دهد که استفاده از اطلاعات چندنما عملکرد مدل را به‌طور قابل‌توجهی بهبود می‌بخشد، اما ترکیب صحیح نماها برای دستیابی به بهترین نتایج ضروری است. در حالی که الگوریتم برای Stack Overflow طراحی شده، قابلیت انتقال به سایر وبسایت‌های پرسش و پاسخ را نیز دارد. این پژوهش محدودیت خاصی را گزارش نکرده و می‌تواند در آینده با استفاده از الگوهای مشابه در سایر پلتفرم‌ها توسعه یابد. روش BoostStack در مقایسه با روش‌های تکنما دقت، بازخوانی و معیار  $F^{12}$  بالاتری را نشان داده است. به طور خاص، الگوریتم BoostStack در آزمایش‌ها بالاترین دقت و F-measure را به‌دست آورده و عملکرد بهتری نسبت به روش‌های پایه مانند Boost-NC، Boost-CCAP و Boost-CCAU داشته است. این نتایج تأیید می‌کنند که ترکیب داده‌های ناهمگن (متن، کد و اطلاعات پس‌زمینه) با استفاده از این روش، پیش‌بینی بهترین پاسخ‌ها را به طور چشمگیری بهبود می‌بخشد. دقت الگوریتم BoostStack برابر با ۷۷.۳٪ و معیار F برابر با ۷۶.۲٪ گزارش شده است. این مقادیر نسبت به سایر روش‌های مورد آزمایش، بالاتر بوده و نشان‌دهنده عملکرد برتر این الگوریتم در پیش‌بینی بهترین پاسخ‌ها در پلتفرم‌هایی مانند Stack Overflow است.

توندوکار و همکاران (۲۰۱۸) [۲۵] از یک الگوریتم یادگیری رتبه‌بندی<sup>۱۴</sup> برای پیش‌بینی بهترین پاسخ‌دهنده به سوالات در انجمن‌های پرسش و پاسخ استفاده کردند. آن‌ها مجموعه جامعی از ویژگی‌ها شامل شباهت متنی، شباهت برچسب‌ها، و ویژگی‌های کاربرمحور نظیر تخصص

روش‌های مختلف ایجاد کرده است. آینده پژوهش می‌تواند بر شناسایی عوامل دیگر مؤثر بر سوگیری رأی‌دهی و ادغام ویژگی‌های محتوایی در مدل‌ها تمرکز کند. روش پیشنهادی کاربردهای عملی در بهبود رتبه‌بندی پاسخ‌ها و کاربران متخصص در سیستم‌های پرسش و پاسخ اجتماعی دارد.

در پژوهشی توسط Elalfy و همکاران (۲۰۱۸) [۲۳]، یک مدل ترکیبی برای پیش‌بینی بهترین پاسخ در جوامع پرسش و پاسخ پیشنهاد شده است. این مدل شامل دو ماژول ویژگی‌های محتوایی و غیرمحتوایی است که ویژگی‌هایی همچون محتوای پاسخ، ارتباط میان پاسخ‌ها، و سطح اعتماد کاربران را تحلیل می‌کند. با استفاده از الگوریتم‌های طبقه‌بندی نظیر جنگل تصادفی و رگرسیون لجستیک، این پژوهش نشان داد که مدل ترکیبی در پیش‌بینی پاسخ‌های باکیفیت عملکرد بهتری نسبت به مدل‌های صرفاً محتوایی دارد. اگرچه این روش دقت بالایی را ارائه می‌دهد، محدودیت‌هایی همچون نیاز به داده‌های گسترده و هزینه پردازشی مطرح است. شکاف پژوهشی در این مطالعه، توجه ناکافی به نوع سوالات به‌عنوان ویژگی تأثیرگذار است که می‌تواند در پژوهش‌های آتی بررسی شود.

Zheng و Li (۲۰۱۷) [۲۴] روشی نوین به نام BoostStack را برای پیش‌بینی بهترین پاسخ در Stack Overflow ارائه کردند. این روش با بهره‌گیری از داده‌های ناهمگن مانند متن، کدهای نمونه و اطلاعات پس‌زمینه کاربران، از یادگیری چندنما استفاده می‌کند. الگوریتم BoostStack با ترکیب Adaboost به‌عنوان طبقه‌بند سطح پایه و یک شبکه عصبی دو لایه به‌عنوان

<sup>۱۴</sup> LambdaMART<sup>۱۵</sup> F-measure

مثال، بازیابی اسناد تنها در ۷۷.۸٪ موارد پاسخ درست را بازیابی می‌کند. از محدودیت‌ها می‌توان به وابستگی به داده‌های ویکی‌پدیا و عدم امکان تعمیم به سایر حوزه‌ها اشاره کرد.

هو و همکاران (Hu et al., ۲۰۲۳) [۲۷] مدلی برای انتخاب پاسخ در سیستم‌های پرسش و پاسخ جامعه‌محور<sup>۲۴</sup> به نام شبکه توجه متقاطع پرسش-پاسخ<sup>۲۵</sup> یا QAN ارائه کردند که از مدل‌های پیش‌آموزش‌یافته<sup>۲۶</sup> و مدل‌های زبانی بزرگ<sup>۲۷</sup> بهره می‌برد. این مدل با استفاده از BERT برای رمزگذاری کلمات و مکانیزم توجه متقاطع<sup>۲۸</sup>، ویژگی‌های تعاملی بین سوالات و پاسخ‌ها را استخراج می‌کند. همچنین، از دانش خارجی تولیدشده توسط مدل‌های زبانی بزرگ برای بهبود دقت انتخاب پاسخ استفاده می‌شود.

مدل QAN بر روی دو مجموعه داده SemEval<sup>۲۰۱۵</sup> و SemEval<sup>۲۰۱۷</sup> به عملکرد پیشرفته‌ای دست یافت و در مقایسه با مدل‌های پایه در سه معیار ارزیابی (F1، دقت Acc و میانگین دقت MAP) بهتر عمل کرد. بهینه‌سازی ورودی‌های اولیه<sup>۲۹</sup> برای LLM در ابعاد مختلف، مانند طول ورودی و موقعیت سوالات و پاسخ‌ها،

و در دسترس بودن را به کار بردند. مدل پیشنهادی در مقایسه با روش‌های موجود، عملکرد بهتری داشت و میانگین دقت رتبه‌بندی بالاتری<sup>۱۵</sup> ارائه داد. همچنین، مدل توانست بهترین پاسخ‌دهندگان را در ۰.۴۵ ثانیه به صورت بلادرنگ پیش‌بینی کند. از محدودیت‌های این پژوهش، عدم استفاده از داده‌های بزرگ‌تر و نبود امکان بازتولید نتایج مطالعات قبلی به دلیل دسترسی‌ناپذیری مجموعه داده‌ها بوده است.

چن و همکاران (Chen et al., ۲۰۱۷) [۲۶] سیستمی برای پاسخ به سوالات در دامنه باز<sup>۱۶</sup> با استفاده از ویکی‌پدیا به عنوان تنها منبع دانش<sup>۱۷</sup> ارائه کردند. این سیستم شامل یک بخش بازیابی اسناد<sup>۱۸</sup> مبتنی بر تطبیق دوگانه<sup>۱۹</sup> و TF-IDF و یک مدل شبکه عصبی چندلایه بازگشتی<sup>۲۰</sup> برای استخراج پاسخ‌ها از پاراگراف‌های ویکی‌پدیا است. مدل پیشنهادی در مجموعه داده SQuAD به نتایج پیشرفته‌ای دست یافت و از روش‌های موجود بهتر عمل کرد. استفاده از یادگیری چندوظیفه‌ای<sup>۲۱</sup> و نظارت دور دست<sup>۲۲</sup> عملکرد سیستم را بهبود بخشید.

با این حال، عملکرد در شرایط واقعی به دلیل چالش‌های خوانش ماشینی در مقیاس بزرگ<sup>۲۳</sup> محدود است. برای

<sup>۲۳</sup> Machine Reading at Scale<sup>۱۵</sup> MRR@N<sup>۲۴</sup> Community Question Answering<sup>۱۶</sup> Open-Domain Question Answering<sup>۲۵</sup> Question-Answer Cross Attention Network<sup>۱۷</sup> Knowledge Source<sup>۲۶</sup> Pre-trained Models<sup>۱۸</sup> Document Retrieval<sup>۲۷</sup> Large Language Models<sup>۱۹</sup> Bigram Matching<sup>۲۸</sup> Cross Attention Mechanism<sup>۲۰</sup> Multi-layer Recurrent Neural Network<sup>۲۹</sup> Prompts<sup>۲۱</sup> Multitask Learning<sup>۲۲</sup> Distant Supervision

دقت انتخاب پاسخ‌ها را بهبود داده است.

با این حال، چالش‌های مربوط به تعامل جداگانه سوالات و پاسخ‌ها در روش‌های سنتی همچنان مطرح است. محدودیت‌های مدل به طور صریح بحث نشده‌اند، اما مطالعات نشان داده‌اند که هر بخش از مدل QAN برای دستیابی به عملکرد مطلوب ضروری است. بهبود سرعت محاسباتی و استفاده از مدل‌های پیش‌آموزش‌یافته جدید از اولویت‌های کارهای آینده معرفی شده است.

کاربردهای عملی این پژوهش شامل بهبود کارایی و دقت سیستم‌های CQA در حوزه‌های مختلف مانند خدمات مشتری، آموزش و بازیابی اطلاعات است.

توبا و همکاران (۲۰۱۴) [۲۸] چارچوبی سلسله‌مراتبی برای پیش‌بینی کیفیت پاسخ در سامانه‌های پرسش و پاسخ اجتماعی ارائه داده‌اند که از ترکیب طبقه‌بندی نوع پرسش و تحلیل کیفیت پاسخ استفاده می‌کند. این روش با تکیه بر ویژگی‌های ذاتی متن، از جمله تعداد کلمات بدون توقف<sup>۲۰</sup>، طول پاسخ و امتیاز خوانایی، دقت پیش‌بینی کیفیت پاسخ را بهبود می‌بخشد. نتایج نشان داد که این چارچوب در دسته‌بندی پاسخ‌های با کیفیت پایین، نسبت به روش‌های مبتنی بر یک طبقه‌بند، عملکرد بهتری دارد. این پژوهش از مدل یادگیری نظارت‌شده برای طبقه‌بندی نوع پرسش و یک اپراتور تصمیم‌گیری برای تحلیل کیفیت پاسخ استفاده می‌کند. مدل‌ها بر اساس ویژگی‌های معنایی، آماری، چگالی، و قطبیت احساسات آموزش داده شدند. همچنین سه روش تجمیع داده‌ها (SUM، MAX و STEP-۲) برای ارزیابی کیفیت پاسخ مقایسه شدند. چارچوب پیشنهادی حدود ۷۰ درصد دقت را روی مجموعه داده CQA به دست آورد و عملکرد بهتری نسبت

به مدل‌های پایه ارائه داد. روش MAX بهترین نتایج را در پیش‌بینی کیفیت پاسخ‌ها نشان داد. همچنین، تحلیل‌ها نشان داد که کیفیت پاسخ‌ها می‌تواند بر اساس طول، ساختار، و محتوای آن‌ها سنجیده شود. مجموعه داده مورد استفاده به طور طبیعی نامتوازن بود (۸۰ درصد پاسخ‌های با کیفیت بالا)، که می‌تواند بر دقت چارچوب اثر بگذارد. همچنین، چالش‌هایی در تشخیص اسم‌ها در پاسخ‌ها وجود داشت که ممکن است فرآیند طبقه‌بندی را تحت تأثیر قرار دهد. پژوهشگران پیشنهاد می‌دهند که از ویژگی‌های اضافی، نظیر اطلاعات کاربر و دانش خارجی، برای بهبود دقت پیش‌بینی استفاده شود. همچنین، چارچوب پیشنهادی می‌تواند در سایر انواع محتواهای تولید شده توسط کاربران، نظیر پست‌های وبلاگ یا توییت‌ها، تعمیم یابد. این چارچوب قابلیت استفاده در بهبود کیفیت پاسخ‌ها در سامانه‌های پرسش و پاسخ نظیر Yahoo! Answers را دارد، که منجر به ارتقاء تجربه کاربر و کیفیت اشتراک‌گذاری دانش می‌شود.

کالفاتو و همکاران (۲۰۱۹، Calefato et al.) [۲۹] به ارزیابی عملکرد الگوریتم‌های مختلف طبقه‌بندی برای پیش‌بینی بهترین پاسخ در سایت Stack Overflow پرداختند. آن‌ها از معیارهای مستقل از آستانه مانند AUC و Balance استفاده کرده و ویژگی‌های کلیدی مؤثر بر دقت پیش‌بینی را شناسایی کردند. نتایج نشان داد که انتخاب الگوریتم طبقه‌بندی و تنظیم خودکار پارامترها تأثیر قابل توجهی بر عملکرد مدل‌ها دارد. در پیش‌بینی درون-پلتفرمی<sup>۲۱</sup>، مدل‌های برتر به عملکردی با میانگین AUC حدود ۰.۹ دست یافتند. با این حال، عملکرد

<sup>۲۱</sup> within-platform

<sup>۲۰</sup> non-stop words

تقویت پرسش<sup>۳۴</sup>، ایجاد برجسب<sup>۳۵</sup>، و توصیه پاسخ<sup>۳۶</sup>. در مرحله اول، پرسش‌های شفاف‌سازی با استفاده از مدل دنباله به دنباله<sup>۳۷</sup> ایجاد می‌شوند. در مرحله دوم، برجسب‌هایی برای داده‌ها شامل مثبت، خنثی، خنثی- و منفی با استفاده از چهار قاعده تجربی ساخته می‌شود. در نهایت، مرحله توصیه پاسخ با استفاده از شبکه عصبی کانولوشنی<sup>۳۸</sup> امتیاز تطابق بین پرسش و پاسخ‌ها را محاسبه می‌کند. نتایج تجربی نشان می‌دهند که این رویکرد به طور چشمگیری از چندین مدل پایه موجود در ارزیابی خودکار پیشی می‌گیرد و در ارزیابی‌های انسانی نیز اثربخش بوده است. اگرچه این روش در حل مشکل تشنگی پاسخ موفقیت‌آمیز است، محدودیت‌های آن شامل تمرکز بر سؤالات مرتبط با زبان‌های برنامه‌نویسی پایتون و جاوا در سایت Stack Overflow به دلیل هزینه بالای مطالعه تجربی است. این رویکرد می‌تواند در سایت‌های پرسش و پاسخ فنی مانند Stack Overflow به توسعه‌دهندگان کمک کند تا پاسخ‌های مرتبط‌تر را شناسایی کنند و زمان جستجوی غیرضروری برای یافتن پاسخ‌های حل نشده را کاهش دهند.

در مطالعه Lou و همکاران (۲۰۲۳) [۳۱]، یک مرور جامع درباره‌ی پژوهش‌های مرتبط با پیروی از دستورات<sup>۳۹</sup> در پردازش زبان طبیعی<sup>۴۰</sup> ارائه شده است. این مطالعه موضوعاتی مانند تعریف دستورات وظیفه، مدل‌سازی دستورات، مجموعه داده‌های رایج و معیارهای ارزیابی و

مدل‌ها در پیش‌بینی بین-پلتفرمی<sup>۳۲</sup> با افتی ۲۰ تا ۲۷ درصدی همراه بود. همچنین تحلیل ویژگی‌ها نشان داد که ویژگی‌های مرتبط با امتیاز، قوی‌ترین پیش‌بین‌ها هستند. محدودیت‌های این پژوهش شامل تکیه بیش از حد مدل‌ها به ویژگی‌های مبتنی بر امتیاز و عدم توجه به تأثیر زمان‌بندی در فرآیند پرسش و پاسخ است. این محدودیت‌ها ممکن است به کاهش عملکرد مدل‌ها در پیش‌بینی بین-پلتفرمی منجر شوند. در نتیجه، نویسندگان توصیه می‌کنند که تحقیقات آینده بر استفاده از ویژگی‌ها و تکنیک‌های جدید تمرکز کنند و از روش‌هایی مانند پردازش زبان طبیعی و تحلیل احساسات برای بهبود پیش‌بینی استفاده شود. این مطالعه کاربردهای عملی مهمی دارد، از جمله کمک به انتقال محتوا از فروم‌های قدیمی به پلتفرم‌های مدرن پرسش و پاسخ و تضمین کیفیت دانش جمع‌سپاری شده در این سایت‌ها. یافته‌ها می‌توانند برای طراحان این پلتفرم‌ها در ساخت ابزارهای خودکار انتخاب بهترین پاسخ مفید باشند.

در پژوهشی که توسط Gao و همکاران (۲۰۲۰) [۳۰] انجام شد، یک رویکرد مبتنی بر یادگیری عمیق به نام DeepAns برای شناسایی پاسخ‌های مرتبط‌تر در سایت‌های پرسش و پاسخ فنی معرفی گردید. این رویکرد به حل مشکل "تشنگی پاسخ"<sup>۳۳</sup> کمک می‌کند و نشان داده شده که عملکرد بهتری نسبت به مدل‌های مرسوم دارد. DeepAns از سه مرحله اصلی تشکیل شده است:

<sup>۳۷</sup> sequence-to-sequence<sup>۳۸</sup> CNN<sup>۳۹</sup> Instruction Following<sup>۴۰</sup> Natural Language Processing<sup>۳۲</sup> cross-platform<sup>۳۳</sup> answer hungriness<sup>۳۴</sup> question boosting<sup>۳۵</sup> label establishment<sup>۳۶</sup> answer recommendation



کرده‌اند. از محدودیت‌های این پژوهش می‌توان به تعصب در انتخاب شرکت‌کنندگان، گزارش‌دهی ذهنی درد و عدم تعمیم‌پذیری به اقلیم‌های مختلف اشاره کرد. این پژوهش نشان می‌دهد که استفاده از داده‌های جمع‌آوری‌شده توسط ابزارهای شهروند-علمی<sup>۴۳</sup> می‌تواند درک ما از ارتباط بین متغیرهای محیطی و درد را بهبود بخشد و به توسعه پیش‌بینی‌های مرتبط با درد کمک کند.

والنتین (۲۰۲۲) [۳۳] از مدل‌های یادگیری عمیق شامل لایه‌های تعبیه<sup>۴۴</sup> و شبکه‌های حافظه طولانی‌مدت<sup>۴۵</sup> برای پیش‌بینی پاسخ‌های انتخاب‌شده توسط کاربران در Stack Overflow استفاده کرده است. این مدل‌ها به همراه ویژگی‌های عددی برای رتبه‌بندی پاسخ‌ها به کار گرفته شدند. هرچند استفاده از این روش‌ها باعث بهبود عملکرد نسبت به برخی مدل‌های پایه شد، اما نتوانستند عملکرد بهترین مدل‌های پایه را بهبود دهند. عملکرد مدل یادگیری عمیق در مقایسه با مدل‌های پایه به دلیل وجود تعداد زیادی کلمات خارج از واژگان<sup>۴۶</sup> یا OOV که ناشی از حضور کدهای برنامه‌نویسی در متن بود، محدود شد. همچنین، محدودیت طول توالی<sup>۴۷</sup> در مرحله تعبیه باعث از دست رفتن بخشی از اطلاعات شده است. نتایج نشان داد که مدل‌های پایه بر اساس ویژگی‌های عددی و متن، عملکرد بهتری نسبت به مدل‌های یادگیری عمیق داشتند. دقت دقیق مدل‌های یادگیری عمیق و پایه در مقاله به طور کمی ارائه نشده است، اما بیان شده که مدل یادگیری عمیق نتوانست به عملکرد بهترین مدل پایه که از

چالش‌های پیش روی این حوزه را بررسی می‌کند. نویسندگان به اهمیت هم‌خوانی با اهداف پیش‌آموزش مدل‌های زبانی بزرگ و افزایش مقیاس مدل برای بهبود عملکرد اشاره کرده‌اند. همچنین، نتایج نشان می‌دهد که حتی مدل‌های پایه نیز می‌توانند عملکرد بهتری نسبت به مدل‌های کوچک‌تر تنظیم‌شده داشته باشند. محدودیت‌های این پژوهش شامل هزینه بالا و تأثیرات زیست‌محیطی مقیاس‌بندی مدل‌ها و ناکارآمدی روش‌های ارزیابی فعلی است. پیشنهاداتی برای آینده شامل توسعه روش‌های نمایشی جدید وظایف و رفع مشکلات مربوط به منفی‌سازی در دستورات ارائه شده است.

در مطالعه Sakai و همکاران (۲۰۱۱) [۳۲]، از معیارهای مرتبط با درجه‌بندی مرتبط بودن<sup>۴۱</sup> برای ارزیابی انتخاب پاسخ‌ها در سیستم‌های پرسش و پاسخ اجتماعی استفاده شده است. این پژوهش، داده‌های روزانه‌ی بیش از ۲۶۵۸ شرکت‌کننده با شرایط درد مزمن را در یک بازه زمانی ۱۵ ماهه تحلیل کرده و روابط معناداری بین درد و متغیرهای جوی مانند رطوبت نسبی، فشار و سرعت باد شناسایی کرده است. نتایج نشان داد که رطوبت نسبی بیشترین ارتباط را با درد دارد (احتمال وقوع رویداد درد با هر افزایش انحراف معیار در رطوبت نسبی ۱۲ درصد افزایش می‌یابد). روش‌های استفاده‌شده شامل طراحی مورد-مقابله‌ای<sup>۴۲</sup> و مدل‌های رگرسیون لجستیک شرطی بوده که تأثیر عواملی مانند خلق، فعالیت فیزیکی، و زمان سپری‌شده در خارج از منزل را نیز در تحلیل‌ها لحاظ

<sup>۴۵</sup> LSTM<sup>۴۶</sup> Out-of-Vocabulary<sup>۴۷</sup> Sequence Length<sup>۴۱</sup> Graded-Relevance Metrics<sup>۴۲</sup> Case-Crossover Design<sup>۴۳</sup> Citizen-Science Experiments<sup>۴۴</sup> Embedding

استفاده از الگوریتم‌های یادگیری ماشین، ویژگی‌هایی شامل متنی و غیرمتنی (مانند تنوع لغوی، شهرت کاربر و شباهت پاسخ‌ها) را بررسی می‌کند و پاسخ‌های باکیفیت را در برگه‌ای به نام "پاسخ‌های امیدبخش"<sup>۵۱</sup> نمایش می‌دهد. مدل پیشنهادی با استفاده از ۲۶ ویژگی مختلف، از جمله تنوع اسمی و صفتی و شباهت بین پاسخ‌ها، دقت پیش‌بینی  $F1\text{-Score}$  ۸۰٪ و بازخوانی ۹۰٪ را ارائه داده است. مدل Gradient Boosting در دسته‌بندی و مدل رگرسیون جنگل تصادفی<sup>۵۲</sup> در پیش‌بینی تعداد رأی‌ها عملکرد بهتری نسبت به مدل‌های بدون طبقه‌بند داشت. این سیستم، با میانگین دقت رتبه‌بندی<sup>۵۳</sup> برابر ۰.۴۸۳ و میانگین دقت متوسط<sup>۵۴</sup> برابر ۰.۴۸۳، توانایی پیش‌بینی بالای کیفیت پاسخ‌ها را نشان داد. محدودیت اصلی پژوهش استفاده از داده‌های محدود و تمرکز بر یک مجموعه داده خاص بود، اما پتانسیل تعمیم به سایر پلتفرم‌های CQA را نیز دارد.

زوراتو و همکاران (۲۰۲۳) [۳۶] تأثیر ویژگی‌های مختلف پاسخ‌ها را در پیش‌بینی بهترین پاسخ‌ها در انجمن‌های پرسش‌وپاسخ بررسی کرده‌اند. آن‌ها از الگوریتم‌های کلاسیک یادگیری ماشین مانند بیز ساده<sup>۵۵</sup>، رگرسیون لجستیک<sup>۵۶</sup> و جنگل تصادفی بر روی مجموعه داده‌ای شامل ۹۴۲۸ پست از Stack Overflow استفاده کردند. نتایج نشان داد که قابلیت خوانایی<sup>۵۷</sup> مهم‌ترین ویژگی برای تشخیص بهترین پاسخ‌ها است و به تنهایی

روش‌های کیسه کلمات<sup>۴۸</sup> استفاده می‌کرد، نزدیک شود. پیشنهاد شده که در تحقیقات آینده از مدل‌های پیشرفته‌تر مانند خانواده مدل‌های BERT استفاده شود و همچنین روش‌هایی برای شناسایی خودکار و حذف کدهای برنامه‌نویسی توسعه یابد تا مشکل OOV حل شود.

روی و همکاران (۲۰۲۳) [۳۴] در مقاله‌ای به بررسی سیستم‌های پرسش و پاسخ اجتماعی با استفاده از یادگیری ماشین و یادگیری عمیق پرداختند. آن‌ها پژوهش‌های موجود را در سه بخش اصلی دسته‌بندی کردند: پرسش‌ها، پاسخ‌ها و کاربران. این مطالعه نشان داد که اکثر تحقیقات بر کیفیت پرسش‌ها، کیفیت پاسخ‌ها و شناسایی کاربران خبره تمرکز دارند. در حالی که یادگیری ماشین در این حوزه غالب است، استفاده از یادگیری عمیق در حال افزایش است. مدل‌های ترکیبی و چندوجهی کمتر استفاده شده‌اند، و بسیاری از مدل‌ها تنها برای یک پلتفرم خاص آزمایش شده‌اند. از جمله نتایج، مدل‌های پیشنهادی به همبستگی پیرسون<sup>۴۹</sup> ۰.۶۲ و کوهن کاپا<sup>۵۰</sup> ۰.۸۶ دست یافتند. پژوهشگران نیاز به تحقیقات بیشتر در زمینه مدل‌های عمومی که قابلیت اجرا بر روی پلتفرم‌های مختلف را داشته باشند و استفاده از تکنیک‌های ترکیبی و چندوجهی را برجسته کردند.

روی و همکاران (۲۰۱۸) [۳۵] سیستمی را برای رتبه‌بندی پاسخ‌ها در سایت‌های پرسش و پاسخ اجتماعی بر اساس کیفیت آن‌ها پیشنهاد کردند. این سیستم، با

<sup>۵۲</sup> MRR<sup>۵۴</sup> MAP<sup>۵۵</sup> Naïve Bayes<sup>۵۶</sup> Logistic Regression<sup>۵۷</sup> Readability<sup>۴۸</sup> Bag of Words<sup>۴۹</sup> Pearson correlation<sup>۵۰</sup> Cohen kappa<sup>۵۱</sup> Promising Answers<sup>۵۲</sup> Random Forest Regression

۱۳.۵٪ موضوعها به عنوان "نامناسب" شناخته شدند که نسبت به ۱۶.۵٪ در LDA بهبود داشت اما از لحاظ آماری معنادار نبود ( $p=0.05$ ). پژوهشگران همچنین به این نتیجه رسیدند که هیچ متدی قادر به پیش‌بینی کامل قضاوت‌های انسانی نیست. بهبود کیفیت معنایی موضوعها، خصوصاً در سطح پایین و مقیاس‌پذیری برای مجموعه داده‌های بسیار بزرگ، به عنوان مسیرهای آتی تحقیق پیشنهاد شد. این پژوهش کاربردهای عملی در حوزه‌هایی مانند بازیابی اطلاعات و خلاصه‌سازی متون دارد و می‌تواند برای ایجاد مدل‌هایی با موضوعات خاص و باکیفیت در مجموعه‌های داده حوزه‌ای استفاده شود.

در مطالعه‌ای که توسط Molino و همکاران (۲۰۱۶) [۳۸] انجام شد، یک رویکرد جدید برای پیش‌بینی بهترین پاسخ در سایت‌های پرسش و پاسخ اجتماعی پیشنهاد شده است. این پژوهش از مجموعه ویژگی‌های گسترده شامل ویژگی‌های متنی، کاربری و شبکه‌ای استفاده کرده و بهبود ۱۱ تا ۲۶ درصدی در معیار دقت در رتبه یک<sup>۶۲</sup> نسبت به روش‌های پیشرفته پیشین نشان داده است. این تحقیق نشان داد که ویژگی‌های معنایی توزیعی<sup>۶۳</sup> می‌توانند جایگزین ویژگی‌های پرهزینه و زمان‌بر شباهت زبانی شوند و قدرت پیش‌بینی بالاتری ارائه دهند. همچنین، مشخص شد که کیفیت متن برای پیش‌بینی بهترین پاسخ در سوالات واقعی و ذهنی مناسب‌تر است، در حالی که ویژگی‌های مربوط به پروفایل کاربر برای سوالات بحث‌برانگیز و نظرسنجی‌ها کاربرد بیشتری دارند. محدودیت‌های این تحقیق مورد بررسی قرار نگرفته است،

توانست دقت TPR برابر با ۰.۷۶۷ را به دست آورد. ترکیب قابلیت خوانایی و شباهت محتوایی<sup>۵۸</sup> نیز عملکرد قابل‌توجهی با TPR برابر ۰.۷۴۱ داشت. با این حال، نقش اجتماعی کاربران و تأثیر شهرت آنها به عنوان شکاف پژوهشی معرفی شده و نیازمند بررسی بیشتر است. این پژوهش نشان می‌دهد که بهبود خوانایی می‌تواند به ارتقاء کیفیت محتوای تولید شده و افزایش اعتبار کاربران کمک کند.

میمو و همکاران (۲۰۱۱) [۳۷] در این پژوهش بهبود کیفیت موضوعها در مدل‌های آماری موضوعی<sup>۵۹</sup> را از طریق معرفی یک معیار جدید برای ارزیابی انسجام موضوعها بررسی کردند. آن‌ها از مدل آماری جایگزین *Polya urn* تعمیم‌یافته استفاده کردند که بهبود قابل‌توجهی در انسجام موضوعها و کاهش موضوعهای کم‌کیفیت نشان داد. این مدل بر اساس اطلاعات هم‌رخدادی واژگان طراحی شده و توانست در مجموعه‌ای از ۳۰۰,۰۰۰ خلاصه مقالات و پروژه‌های تحقیقاتی مؤسسات ملی بهداشت<sup>۶۰</sup> عملکرد بهتری نسبت به روش‌های سنتی مانند تخصیص دیریکله نهان ارائه دهد. نتایج نشان داد که این مدل انسجام متوسط موضوعها را بهبود بخشیده و امتیاز انسجام برای ده موضوع با پایین‌ترین امتیاز را افزایش داده است. معیار انسجام پیشنهادی، موضوعهای باکیفیت پایین را بهتر از معیارهایی مانند اندازه موضوع و اطلاعات متقابل نقطه‌ای<sup>۶۱</sup> شناسایی کرد. با این حال، مدل در کاهش کلی تعداد موضوعهای نامناسب موفقیت چندانی نداشت و تنها

<sup>۶۱</sup> PMI<sup>۵۸</sup> Content Similarity<sup>۶۲</sup> P@1<sup>۵۹</sup> Topic Models<sup>۶۳</sup> Distributional Semantics<sup>۶۰</sup> NIH

- اما پیشنهاد شده که تحقیقات آینده به بررسی کارایی این ویژگی‌ها در سایر برنامه‌ها بپردازند.
  - دقت ۸۴.۷۲٪ با Logistic Regression برای پیش‌بینی کیفیت پاسخ.
  - دقت ۸۸٪ با مدل هیبریدی شامل Random Forest و Logistic Regression برای پیش‌بینی بهترین پاسخ.
  - دقت ۹۵.۲٪ با Light Gradient Boosting Machine برای پیش‌بینی بهترین پاسخ.
- در مقایسه با این پژوهش‌ها، مدل پیشنهادی ما با دقت ۹۹.۰۳٪ نه تنها دقت بالاتری ارائه می‌دهد، بلکه از روش‌های پیشرفته‌ای مانند پشته‌ای وزن‌دار و ترکیب ویژگی‌های متنی و غیرمتنی بهره می‌برد. علاوه بر این، استفاده از مدل‌های پیشرفته‌ای مانند XLNet و RoBERTa، امکان استخراج معنای عمیق‌تر از متن و بهبود پیش‌بینی‌ها را فراهم کرده است.

این مدل در مقایسه با پژوهش‌های پیشین از دقت و قابلیت تعمیم‌دهی بالاتری برخوردار باشد. این دستاورد می‌تواند به عنوان یک الگوی جدید در توسعه سیستم‌های پرسش و پاسخ مورد استفاده قرار گیرد و تجربه کاربران در این حوزه را بهبود بخشد.

### ۳-۲- تفاوت و بهبودهای مدل ما نسبت به پژوهش‌های پیشین

در حالیکه بسیاری از پژوهش‌های پیشین در این حوزه مانند Burel و همکاران (۲۰۱۲) [۲۰] بر ویژگی‌های زبانی سطحی یا ویژگی‌های ساده متکی بودند، مدل ما از

انبارکی و جوکار (۲۰۲۱) [۳۹] به ارزیابی و پیش‌بینی کیفیت پاسخ‌ها در شبکه اجتماعی پژوهش‌محور<sup>۶۴</sup> در حوزه مدیریت دانش پرداختند. این مطالعه با تحلیل توصیفی و استنباطی بر روی ۵۴ سوال و ۴۴۳ پاسخ، که توسط ۳۰ متخصص از دو دانشگاه با استفاده از ۱۳ معیار ارزیابی شدند، انجام شد. نتایج نشان داد که متغیرهای ارتباط<sup>۶۵</sup>، کفایت<sup>۶۶</sup> و هماهنگی<sup>۶۷</sup> با ضرایب شانس به ترتیب ۳.۶۲۶، ۳.۴۴۰ و ۳.۱۴۸، بیشترین قدرت را در پیش‌بینی پاسخ‌های صحیح یا نادرست دارند. محدودیت‌های این پژوهش شامل حجم نمونه کوچک و ارزیابی ذهنی توسط متخصصین بود. پیشنهاد شده است که پژوهش‌های آینده بر توسعه روش‌های خودکار و عینی برای ارزیابی کیفیت پاسخ‌ها تمرکز کنند.

### ۳-۱- مقایسه با پژوهش‌های پیشین

پژوهش‌های پیشین در این حوزه از مدل‌های مختلفی مانند SVM، Adaboost، Logistic Regression و شبکه‌های عصبی استفاده کرده‌اند. اما این پژوهش‌ها با محدودیت‌هایی از جمله دقت پایین، عدم توانایی ترکیب ویژگی‌های متنی و غیرمتنی و عدم استفاده بهینه از داده‌های محدود مواجه بوده‌اند. برخی از نتایج پژوهش‌های پیشین عبارت‌اند از:

- دقت ۶۶٪ با استفاده از Adaboost برای پیش‌بینی کیفیت سؤال.

<sup>۶۴</sup> Adequacy<sup>۶۴</sup> ResearchGate<sup>۶۷</sup> Concordance<sup>۶۵</sup> Relevance

اما دقت عددی ارائه نشده است. روش‌های ساده‌تر مانند ACQUA (Gkotsis و همکاران) دقت ۸۴٪ و نرخ بازیابی ۷۰٪ را ثبت کرده‌اند.

در حالی که بسیاری از مدل‌ها مانند روش پیشنهادی توسط Burel و همکاران (۲۰۱۲) [۲۰] محدود به داده‌های خاصی هستند، مدل ما با ترکیب ویژگی‌های متنوع توانایی تعمیم‌پذیری بیشتری دارد.

مدل ما با استفاده از ترکیب یادگیرنده‌ها و وزن‌دهی، از مشکلاتی مانند بیش‌برازش جلوگیری کرده است. این امر در مقایسه با مدل‌های پیچیده‌تر مانند DeepAns (Gao و همکاران) [۳۰] که به دلیل نیاز به هزینه پردازشی بالا محدود شده بودند، یک مزیت محسوب می‌شود.

#### ۴-۳- شکاف‌های پژوهشی که مدل ما آن‌ها را پر کرده است

بسیاری از پژوهش‌ها بر یک نوع داده (متنی یا غیرمتنی) تمرکز داشته‌اند. مدل ما با ترکیب هر دو نوع داده، توانسته است شکاف موجود در تحلیل جامع کیفیت پاسخ‌ها را پر کند.

برخلاف مطالعاتی که از مدل‌های قدیمی‌تر یا ویژگی‌های ساده استفاده کرده‌اند (مانند ویژگی‌های زبانی سطحی در Gkotsis و همکاران، ۲۰۱۵ [۱۶])، مدل ما از جدیدترین تکنیک‌های پردازش زبان طبیعی بهره می‌برد.

با وجود پیچیدگی مدل‌های یادگیری عمیق، مدل ما با استفاده از فرایادگیر مبتنی بر درخت تصمیم و ترکیب هوشمندانه خروجی‌ها، قابلیت تفسیرپذیری را حفظ کرده است. همچنین، برخلاف پژوهش‌هایی که تنها در یک پلتفرم خاص آزمایش شده‌اند، مدل ما پتانسیل کاربرد در حوزه‌های مختلف را دارد.

استفاده از اعتبارسنجی متقاطع در مدل ما برای مقابله با کمبود داده‌های برچسب‌خورده، یک مزیت قابل توجه در

مدل‌های پیشرفته زبانی بهره گرفته که توانایی استخراج معنا و زمینه‌های پیچیده از متن را دارد. این امر باعث می‌شود عملکرد مدل در تحلیل دقیق‌تر پرسش‌ها و پاسخ‌ها بهبود یابد.

برخلاف روش‌هایی که به یک نوع ویژگی متکی بوده‌اند، مانند Zhu و همکاران (۲۰۰۹) [۱۷] (تمرکز بر ویژگی‌های متنی) یا Hadfi و همکاران (۲۰۲۲) [۱۵] (تمرکز بر ویژگی‌های غیرمتنی)، مدل ما ترکیبی از هر دو نوع ویژگی متنی و غیرمتنی را برای افزایش دقت پیش‌بینی استفاده کرده است.

مدل ما با ترکیب چندین یادگیرنده پایه قدرتمند مانند جنگل‌های تصادفی، XLNet و RoBERTa و به‌کارگیری یک فرایادگیر مبتنی بر درخت تصمیم، به دقت بالای ۹۹.۰۳٪ دست یافته است. این ترکیب، در مقایسه با مدل‌های منفرد مانند BoostStack (Zheng و Li، ۲۰۱۷) [۲۴] و یا DFFN (Suggu و همکاران، ۲۰۱۶) [۱۴] انعطاف‌پذیری و دقت بالاتری را تضمین می‌کند.

استفاده از وزن‌دهی مبتنی بر عملکرد هر مدل پایه، نوآوری دیگری است که در مطالعات پیشین کمتر دیده شده است. این تکنیک تأثیر مثبت بر کاهش خطای کلی مدل دارد.

با به‌کارگیری اعتبارسنجی متقاطع و ترکیب مدل‌های پایه، مدل ما توانسته است در شرایط کمبود داده‌های برچسب‌خورده نیز عملکرد مناسبی ارائه دهد، در حالی که برخی پژوهش‌ها مانند Elalfy و همکاران [۲۳] (۲۰۱۸) به محدودیت داده‌های گسترده اشاره کرده‌اند.

#### ۳-۳- دقت و کارایی مدل ما نسبت به پژوهش‌های پیشین

دقت مدل ما بسیار بالاتر از مطالعات پیشین است. به‌عنوان مثال مدل BoostStack دقت ۷۷.۳٪ را گزارش داده است و DFFN عملکرد "state-of-the-art" داشته،

## ۴-۱- ساختار کلی مدل

مدل پیشنهادی در چندین فاز کلیدی طراحی و پیاده‌سازی شده است.

مقایسه با مدل‌هایی است که نیاز به داده‌های گسترده دارند، مانند روش Elalfy و همکاران (۲۰۱۸) [۲۳].

## ۴- مدل پیشنهادی

## ۴-۱-۱- پیش‌پردازش داده

این بخش شامل آماده‌سازی و تمیزسازی داده‌ها برای ورود به مدل می‌شود. در این پژوهش، هدف اصلی ارزیابی کیفیت پاسخ‌ها و انتخاب بهترین پاسخ در یک پلتفرم پرسش و پاسخ آنلاین است. برای تحقق این هدف، ابتدا داده‌های خام جمع‌آوری شده از پلتفرم با دقت مورد بررسی و پردازش قرار گرفتند. در مرحله اول، داده‌ها پاک‌سازی و ساختاردهی شدند تا برای استفاده در فرآیند مدل‌سازی آماده شوند. در ادامه، با تحلیل دقیق ویژگی‌های متن و غیرمتنی مرتبط با سوالات و پاسخ‌ها، مجموعه‌ای از ستون‌های مؤثر و مرتبط برای مدل‌سازی انتخاب شد.

ویژگی‌های انتخاب‌شده شامل اطلاعات متنوعی نظیر متن سوالات و پاسخ‌ها، امتیازات کاربران، تعداد بازدیدها، تعاملات کاربران (مانند رأی‌ها و نظرات)، تغییرات محتوایی و زمانی، و شاخص‌های کیفیت پاسخ‌ها است. این ویژگی‌ها به گونه‌ای انتخاب شده‌اند که نه تنها نمایانگر رفتار کاربران و کیفیت محتوای تولیدشده باشند، بلکه امکان پیش‌بینی دقیق کیفیت پاسخ‌ها را نیز فراهم کنند.

در جدول ۱ ویژگی‌های استخراج‌شده از مجموعه داده نشان داده شده است.

جدول ۱: ویژگی‌های استخراج‌شده از مجموعه داده

نام ویژگی غیر متنی	نام ویژگی متنی
-----------------------	-------------------

مدل پیشنهادی این پژوهش، رویکردی ترکیبی و نظارت‌شده است که از قدرت یادگیری ژرف برای پردازش داده‌های پیچیده بهره‌برداری کرده و در عین حال، با انتخاب ویژگی‌های دستی برای داده‌های غیرمتنی، دقت و کارایی مدل را بهبود می‌بخشد. این ترکیب، که شامل استفاده از مدل‌های قدرتمند زبانی (مانند XLNet و RoBERTa) و روش‌های کلاسیک (مانند جنگل‌های تصادفی) است، توانسته است نیازهای متنوع و پیچیده جوامع پرسش و پاسخ آنلاین را به‌طور مؤثری برآورده کند.

مدل پیشنهادی بر اساس دو اصل کلیدی طراحی شده است: الف) استفاده از یادگیری ژرف برای داده‌های متنی؛ مدل‌های یادگیری ژرف توانایی بالایی در استخراج خودکار ویژگی‌های پیچیده از داده‌های حجیم دارند، به‌ویژه در تحلیل داده‌های متنی که شامل ساختارهای نحوی و معنایی عمیق هستند. ب) انتخاب ویژگی‌های دستی برای داده‌های غیرمتنی؛ در مواردی که داده‌ها ساده‌تر و ساختارمندتر هستند (مانند آمار رأی‌گیری یا مشخصات کاربران)، استفاده از ویژگی‌های دستی می‌تواند ضمن کاهش پیچیدگی، تفسیرپذیری مدل را نیز افزایش دهد.

این مدل برای رفع شکاف‌های موجود در روش‌های پیشین و پاسخ به چالش‌های کلیدی این حوزه، نظیر تنوع در داده‌های متنی و غیرمتنی و نیاز به تفسیرپذیری نتایج، طراحی شده است. با ترکیب این دو رویکرد، مدل پیشنهادی توانسته است دقت پیش‌بینی را به ۹۹.۰۳٪ برساند و عملکردی فراتر از مدل‌های منفرد یا روش‌های ساده‌تر ارائه دهد.

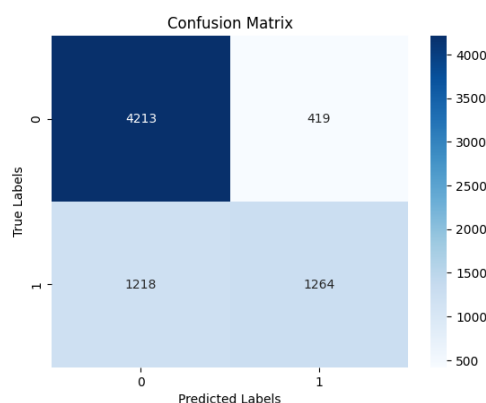
تغییر فاصله زمانی بین درج سوال تا درج پاسخ برحسب روز	
تعداد دفعات عنوان پاسخ	
تعداد دفعات تغییر محتوای پاسخ	
تعداد دفعات تغییر برچسب پاسخ	
پاسخ بسته شده	
پاسخ بازیابی شده	
در نظر گرفته شدن پاسخ به عنوان پاسخ جامعه	
پاسخ حذف شده	
میزان کیفیت پاسخ	

تعداد پاسخها	متن کامل سوال
امتیاز سوال	متن کامل پاسخ
تعداد بازدید برای جفت پست پرسش و پاسخ	
امتیاز پاسخ	
تعداد نظرات درج شده برای پاسخ	
رای به عنوان بهترین پاسخ توسط دیگر کاربران	
تعداد رای مثبت به پاسخ	
تعداد رای منفی به پاسخ	
تعداد رای به بسته شدن پاسخ	
تعداد رای به بازگشایی پاسخ	
تعداد رای به حذف پاسخ	
تعداد رای به بازگرداندن پاسخ حذف شده	

مجموعه داده نهایی پس از آماده سازی، به نسبت ۲۰/۸۰ برای آموزش و آزمون تقسیم شد. همچنین، داده های آموزشی نیز به نسبت ۲۰/۸۰ برای آموزش و اعتبارسنجی مدل تخصیص داده شدند. این تقسیم بندی تضمین می کند که مدل های یادگیری ماشین بتوانند به طور مؤثر آموزش ببینند و عملکرد آن ها به دقت ارزیابی شود.

امتیاز F1	بازخوانی	دقت	
۰.۸۴	۰.۹۱	۰.۷۸	۰
۰.۶۱	۰.۵۱	۰.۷۵	۱
۰.۷۷			دقت

ماتریس درهمی شکل ۱ نشان می‌دهد که مدل XLNet در تشخیص پاسخ‌های غیر برتر (کلاس ۰) عملکرد بهتری داشته است. از مجموع ۴۶۳۲ پاسخ غیر برتر، ۴۲۱۳ پاسخ به درستی تشخیص داده شده‌اند، اما ۱۲۱۸ پاسخ برتر به نادرست به عنوان پاسخ غیر برتر تشخیص داده شده‌اند. این نشان‌دهنده یک عدم تعادل در پیش‌بینی مدل است که باید در آینده با تنظیمات بیشتر بهبود یابد.



شکل ۱: ماتریس درهمی<sup>۶۹</sup> مدل XLNet بر روی مجموعه داده آزمون

RoBERTa نسخه بهینه شده BERT است که با افزایش مقیاس داده‌های آموزشی و بهینه‌سازی تنظیمات، توانسته عملکرد خود را در وظایف پردازش زبان طبیعی بهبود

این فرآیند پیش‌پردازش و انتخاب ویژگی‌ها، داده‌های اولیه را به مجموعه‌ای دقیق، تمیز و غنی تبدیل کرده که می‌تواند پایه‌ای مناسب برای توسعه مدل‌های پیش‌بینی کیفیت پاسخ‌ها باشد. این مرحله، نقشی حیاتی در بهبود عملکرد مدل‌ها و افزایش دقت پیش‌بینی‌ها ایفا می‌کند.

## ۲-۱-۴- ویژگی‌های متنی

در این پژوهش، از مدل‌های پیشرفته زبانی XLNet و RoBERTa برای تحلیل ویژگی‌های متنی پرسش‌ها و پاسخ‌ها در مجموعه داده StackOverflow استفاده شده است. این مدل‌ها که بر پایه معماری تبدیل‌کننده طراحی شده‌اند، قادر به درک روابط پیچیده میان کلمات و عبارات هستند و می‌توانند اطلاعات متنی را با دقت بسیار بالا پردازش و تحلیل کنند.

XLNet یکی از پیشرفته‌ترین مدل‌های زبانی است که با ترکیب مدل‌سازی زبان جایگشت و تکنیک‌های مدرن مانند Transformer-XL، محدودیت‌های مدل‌هایی مانند BERT را رفع کرده و عملکرد برتری در وظایف مختلف پردازش زبان طبیعی ارائه می‌دهد. این مدل با بهره‌گیری از احتمال شرطی توکن‌ها و ترتیب‌های مختلف، توانایی درک دقیق‌تر جملات و متون پیچیده را دارد.

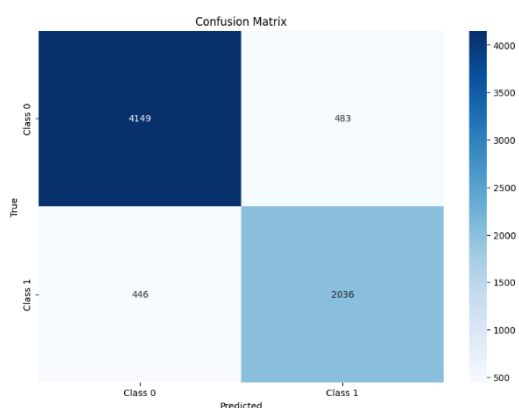
جدول ۲ نشان می‌دهد که دقت نهایی<sup>۶۸</sup> مدل برابر با ۰.۷۷ است، که نشان می‌دهد مدل در ۷۷ درصد موارد به درستی بهترین پاسخ را پیش‌بینی کرده است. با این حال، تفاوت قابل توجهی بین معیارهای Precision و Recall برای کلاس‌های مختلف مشاهده می‌شود.

جدول ۲: گزارش طبقه‌بندی مدل XLNet بر روی مجموعه داده آزمون

<sup>۶۹</sup> Confusion Matrix

<sup>۶۸</sup> Accuracy





شکل ۲ - ماتریس درهمی مدل RoBERTa بر روی مجموعه داده آزمون

هر دو مدل XLNet و RoBERTa توانستند اطلاعات متنی پرسش‌ها و پاسخ‌ها را با دقت بالایی تحلیل کنند. با این حال، RoBERTa با دقت ۸۶.۹۴ درصد عملکرد بهتری نسبت به XLNet (۷۷ درصد) داشت. همچنین، RoBERTa توانست توازن بهتری میان Precision و Recall برای هر دو کلاس برقرار کند.

استفاده از این مدل‌ها به عنوان بخشی از فرآیند کلی شناسایی بهترین پاسخ‌ها، نشان داد که انتخاب مدل مناسب و تنظیم دقیق آن می‌تواند تأثیر چشمگیری بر دقت پیش‌بینی‌ها و کارایی سیستم داشته باشد.

### ۳-۱-۴- ویژگی‌های غیرمتنی

در این بخش، برای شناسایی بهترین پاسخ‌ها در سیستم‌های پرسش و پاسخ، ویژگی‌های غیرمتنی استخراج و مورد تحلیل قرار گرفتند. این ویژگی‌ها که شامل اطلاعات آماری و تعاملی هستند، به کمک مدل جنگل‌های تصادفی شناسایی و برای ورود به مدل‌های ترکیبی پردازش شدند. ویژگی‌های غیرمتنی شامل عواملی می‌شوند که از محتوای متن فراتر رفته و به جنبه‌های آماری، رفتاری و تعاملی کاربران و پاسخ‌ها می‌پردازند. به کارگیری این ویژگی‌ها در کنار ویژگی‌های متنی، به ایجاد مدلی جامع‌تر و با عملکرد بهتر منجر می‌شود. مدل جنگل‌های تصادفی که به صورت خاص برای تحلیل ویژگی‌های غیرمتنی طراحی شده است، به عنوان یکی از اجزای کلیدی مدل ترکیبی به کار گرفته شد. در ادامه،

بخشد. این مدل با تمرکز بر استخراج ویژگی‌های متنی دقیق‌تر، به تحلیل و تفسیر بهتر داده‌های متنی کمک می‌کند.

نتایج نشان داد که RoBERTa عملکرد بهتری نسبت به XLNet دارد و به دقت نهایی ۸۶.۹۴ درصد دست یافته است. این مدل توانایی بالایی در تشخیص پاسخ‌های برتر و غیر برتر داشته و توازن مناسبی میان معیارهای دقت، بازیابی و  $F_1$  برقرار کرده است.

گزارش طبقه‌بندی نیز شامل معیارهایی نظیر دقت، فراخوانی، امتیاز  $F_1$  و دقت کلی است. این گزارش به خوبی نقاط قوت و ضعف مدل در پیش‌بینی هر یک از کلاس‌ها را مشخص می‌کند. در جدول ۳ گزارش طبقه‌بندی مدل RoBERTa را بر روی مجموعه داده آزمون مشاهده می‌کنیم.

جدول ۳: گزارش طبقه‌بندی مدل RoBERTa بر روی مجموعه داده آزمون

	دقت	بازخوانی	امتیاز $F_1$
۰	۰.۹۰	۰.۹۰	۰.۹۰
۱	۰.۸۱	۰.۸۲	۰.۸۱
دقت			۰.۸۷

دقت نهایی مدل RoBERTa بر روی مجموعه داده آزمون برابر با ۸۶.۹۴ درصد به دست آمد که نشان‌دهنده کارایی مناسب این مدل در مقایسه با سایر روش‌های مورد استفاده در این پروژه بود.

در شکل ۲ ماتریس درهمی مدل RoBERTa را بر روی مجموعه داده آزمون مشاهده می‌کنیم.

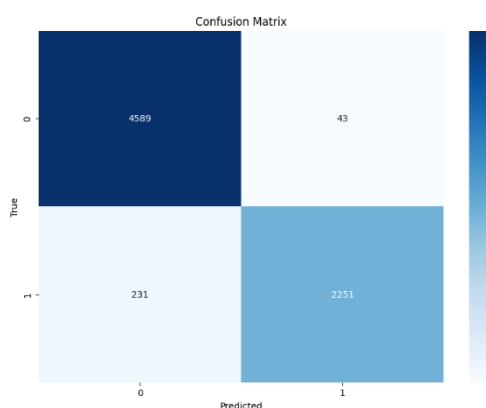
مانند بسته شدن، باز شدن مجدد، حذف شدن و بازگردانی پاسخ و یا حتی تغییر متن و عنوان. همچنین تغییرات در برچسبها نیز جز این عملیات به شمار می رود	فعالیت های مدیریتی
---	--------------------

مهم ترین ویژگی های غیرمتنی که در این پژوهش مورد استفاده قرار گرفتند در جدول ۴ نشان داده شده است.

جدول ۴: ویژگی های غیرمتنی مهم مورد استفاده قرار گرفته در مدل جنگل های تصادفی

ویژگی غیرمتنی	توضیح
امتیاز پاسخ	امتیازی که توسط کاربران به پاسخ داده شده است
امتیاز سوال	امتیازی که توسط کاربران به سوال داده شده است
تعداد بازدها	تعداد دفعاتی که سوال مشاهده شده است.
تعداد پاسخ ها	تعداد پاسخ هایی که به یک سوال داده شده است
تعداد نظرات پاسخ	تعداد نظراتی که برای یک پاسخ داده شده است
تاریخ ایجاد پاسخ	زمان ایجاد پاسخ و مدت زمانی که از انتشار آن گذشته است
فعالیت های امتیازی	مانند تعداد رای های مثبت و منفی که پاسخ دریافت کرده است و یا رأی به بسته شدن یا باز شدن و حتی رأی به حذف شدن یا نشدن

پس از آموزش و ارزیابی مدل، نتایج نشان دادند که مدل جنگل های تصادفی توانایی بسیار بالایی در تشخیص ویژگی های غیرمتنی مرتبط با کیفیت پاسخ ها دارد. دقت مدل، که در جدول ۴ گزارش شده است، برای مجموعه داده آزمون به ۹۶ درصد رسید، که نشان دهنده عملکرد قوی مدل در شناسایی بهترین پاسخ ها است. مدل توانست ۴۵۸۹ نمونه غیر برتر را به درستی شناسایی کند، در حالی که تنها ۴۳ نمونه را به اشتباه به عنوان برتر طبقه بندی کرد. با این حال، ضعف مدل در کاهش نرخ False Negative همچنان قابل توجه است. مقادیر دقیق ماتریس درهمی در شکل ۳ نشان داده شده است.



شکل ۳: ماتریس درهمی مدل جنگل های تصادفی بر روی مجموعه داده آزمون

روی داده‌های آزمون ارائه دهد.

جدول ۴: گزارش طبقه‌بندی مدل جنگل‌های تصادفی بر روی

مجموعه داده آزمون

امتیاز F <sub>1</sub>	بازخوانی	دقت	
۰.۹۷	۰.۹۹	۰.۹۵	۰
۰.۹۴	۰.۹۱	۰.۹۸	۱
۰.۹۶		دقت	

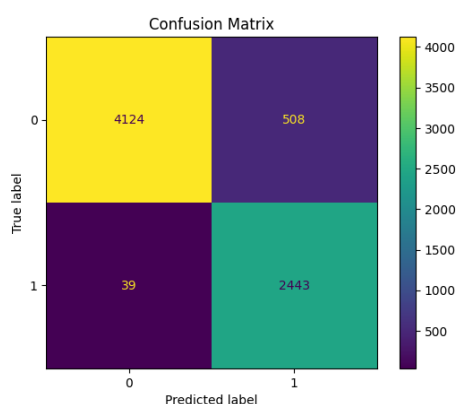
ماتریس درهمی شکل ۴ و گزارش طبقه‌بندی جدول ۵ بیانگر این موضوع است. ماتریس درهمی شکل ۴ نشان می‌دهد که مدل توانسته است تعداد ۴۱۲۴ نمونه از کلاس ۰ و ۲۴۴۳ نمونه از کلاس ۱ را به درستی دسته‌بندی کند. با این حال، تعداد ۵۰۸ نمونه از کلاس ۰ به نادرست به کلاس ۱ و ۳۹ نمونه از کلاس ۱ به نادرست به کلاس ۰ طبقه‌بندی شده‌اند.

جدول ۵ نشان می‌دهد که دقت کلی مدل برابر با ۹۲.۳۱٪

است. همچنین، معیارهای دقت، بازخوانی و امتیاز F<sub>1</sub> برای کلاس‌های ۰ و ۱ محاسبه شده‌اند.

جدول ۵: گزارش طبقه‌بندی مدل رأی‌گیری اکثریت با مدل‌های پایه XLNet، جنگل‌های تصادفی و RoBERTa بر روی مجموعه داده آزمون

امتیاز F <sub>1</sub>	بازخوانی	دقت	
۰.۹۴	۰.۸۹	۰.۹۹	۰
۰.۹۰	۰.۹۸	۰.۸۳	۱
۰.۹۲		دقت	



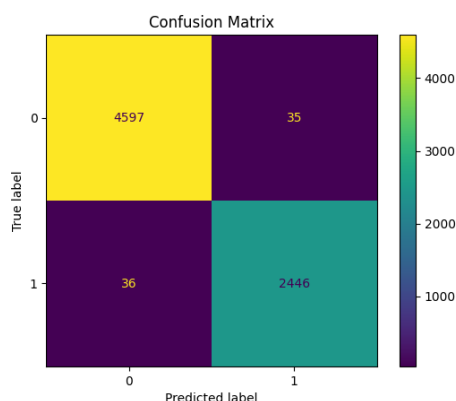
این نتایج نشان می‌دهد که مدل جنگل‌های تصادفی، با استفاده از ویژگی‌های غیرمتنی، به خوبی قادر به تفکیک پاسخ‌های با کیفیت بالا و پایین است و نقش کلیدی در تقویت مدل ترکیبی دارد.

#### ۴-۱-۴- ترکیب مدل‌ها: به کارگیری رأی‌گیری اکثریت با مدل‌های پایه

در این بخش، تکنیک رأی‌گیری اکثریت به عنوان یک روش یادگیری جمعی برای ترکیب پیش‌بینی‌های مدل‌های مختلف به کار گرفته شده است. هدف اصلی این روش، بهبود عملکرد کلی و افزایش دقت پیش‌بینی‌ها از طریق ترکیب خروجی مدل‌های پایه است. مدل‌های به کاررفته شامل XLNet و RoBERTa برای تحلیل ویژگی‌های متنی و جنگل‌های تصادفی برای تحلیل ویژگی‌های غیرمتنی هستند. در این پژوهش، از روش رأی‌گیری سخت استفاده شده است که بر اساس بیشترین تعداد آرا، کلاس نهایی را تعیین می‌کند.

برای ارزیابی مدل رأی‌گیری اکثریت، پیش‌بینی‌های خروجی مدل‌های پایه XLNet، RoBERTa و جنگل‌های تصادفی بر روی داده‌های آموزشی، اعتبارسنجی و آزمون ترکیب شدند. این روش از نقاط قوت هر مدل بهره‌برداری کرده و نقاط ضعف آن‌ها را پوشش داده است. حال نتایج ارزیابی نشان می‌دهند که روش رأی‌گیری اکثریت توانسته است عملکرد قابل توجهی بر

امتیاز  $F_1$  برای هر کلاس در جدول ۶ ارائه شده‌اند و بیانگر عملکرد فوق‌العاده مدل در شناسایی صحیح نمونه‌ها هستند. در شکل ۵، تعداد بسیار کمی از نمونه‌ها به‌نادرست پیش‌بینی شده‌اند که این امر نشان‌دهنده کاهش خطاهای مدل است.



شکل ۵: ماتریس درهمی مدل یادگیری جمعی پشته‌ای وزن‌دار با مدل‌های پایه XLNet، جنگل‌های تصادفی، RoBERTa و رأی‌گیری اکثریت و مدل فرایادگیر درخت تصمیم بر روی مجموعه داده آزمون

جدول ۶: گزارش طبقه‌بندی مدل یادگیری جمعی پشته‌ای وزن‌دار با مدل‌های پایه XLNet، جنگل‌های تصادفی، RoBERTa و رأی‌گیری اکثریت و مدل فرایادگیر درخت تصمیم بر روی مجموعه داده آزمون

امتیاز $F_1$	بازخوانی	دقت	
۰.۹۹	۰.۹۹	۰.۹۹	۰
۰.۹۹	۰.۹۹	۰.۹۹	۱

شکل ۴: ماتریس درهمی مدل رأی‌گیری اکثریت با مدل‌های پایه XLNet، جنگل‌های تصادفی و RoBERTa بر روی مجموعه داده آزمون

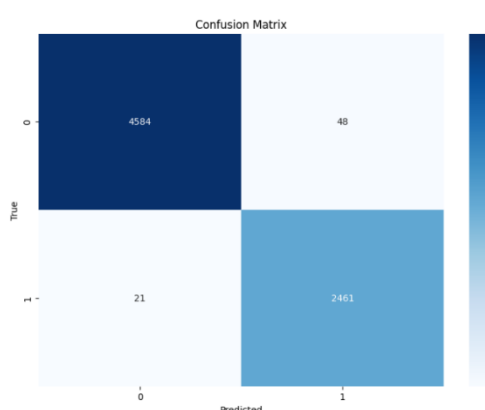
مدل رأی‌گیری اکثریت، با ترکیب پیش‌بینی‌های سه مدل پایه XLNet، RoBERTa و جنگل‌های تصادفی، توانسته است بهبود چشمگیری در دقت و تعمیم‌پذیری پیش‌بینی‌ها ایجاد کند. این روش با ترکیب قابلیت‌های متمایز هر مدل، عملکردی قوی و پایدار ارائه داده است.

#### ۵-۱-۴- فرایادگیر: به‌کارگیری مدل یادگیری جمعی پشته‌ای وزن‌دار با مدل‌های پایه و فرایادگیر درخت تصمیم

در این بخش، از تکنیک یادگیری جمعی پشته‌ای برای ترکیب مدل‌های مختلف به‌منظور بهبود دقت پیش‌بینی‌ها و مدیریت نتایج ترکیبی استفاده شده است. پشته‌ای یک روش پیشرفته در یادگیری ماشین است که در آن چندین مدل پایه برای یادگیری ویژگی‌های متنوع داده‌ها به کار گرفته می‌شوند. پیش‌بینی‌های این مدل‌ها به‌عنوان ورودی به یک مدل فرایادگیر<sup>۷۰</sup> داده می‌شوند تا پیش‌بینی نهایی انجام شود. در این پژوهش، مدل‌های پایه شامل XLNet و RoBERTa برای تحلیل و استخراج ویژگی‌های معنایی متنی با استفاده از یادگیری ژرف، جنگل‌های تصادفی، برای شناسایی و تحلیل ویژگی‌های غیرمتنی و رأی‌گیری اکثریت، به‌عنوان یک ویژگی جدید، ترکیب پیش‌بینی‌های مدل‌های پایه، می‌شود. همچنین از درخت تصمیم به‌عنوان مدل فرایادگیر استفاده شد چراکه درخت تصمیم با توانایی پردازش ترکیب پیچیده ویژگی‌ها را دارد. با توجه به جدول ۶، مدل پشته‌ای وزن‌دار به دقت ۹۹٪ بر روی داده‌های آزمون دست یافت که نشان‌دهنده پیش‌بینی‌های بسیار دقیق و قابل‌اعتماد است. معیارهای دقت، بازخوانی و

بهترین پاسخ انتخاب شود و از مشکلات پیشین جلوگیری شود.

پس از اعمال این رویکرد، دقت و کارایی مدل به طور چشمگیری بهبود یافت. خروجی نهایی مدل، شامل شاخص‌های عملکرد و ماتریس درهمی، نشان داد که مدل توانسته است با دقت ۹۹.۰۳٪ عملکردی بسیار قوی ارائه دهد. نتایج دقیق در جدول ۷ و ماتریس درهمی شکل ۶ به طور کامل ارائه شده است.



شکل ۶: ماتریس درهمی مدل یادگیری جمعی پشته‌ای وزن‌دار با مدل‌های پایه XLNet، جنگل‌های تصادفی، RoBERTa و رأی‌گیری اکثریت و مدل فرایادگیر درخت تصمیم بر روی مجموعه داده آزمون

جدول ۷: گزارش طبقه‌بندی مدل یادگیری جمعی پشته‌ای وزن‌دار با مدل‌های پایه XLNet، جنگل‌های تصادفی، RoBERTa و رأی‌گیری اکثریت و مدل فرایادگیر درخت تصمیم و استفاده از خروجی درجه وقوع هر کلاس به جای پیش‌بینی دقیق کلاس در مدل فرایادگیر بر روی مجموعه داده آزمون

امتیاز F1	بازخوانی	دقت	
۰.۹۹	۱.۰۰	۰.۹۹	۰

۰.۹۹			دقت
------	--	--	-----

روش پشته‌ای با ترکیب مدل‌های قدرتمند XLNet، RoBERTa، جنگل‌های تصادفی و رأی‌گیری اکثریت توانست نقاط قوت هر مدل را بهینه کرده و نقاط ضعف آن‌ها را پوشش دهد. استفاده از درخت تصمیم به عنوان مدل فرایادگیر موجب بهبود قابل توجه دقت پیش‌بینی‌ها شد.

این مدل، با دستیابی به دقت ۹۹٪، نشان داد که ترکیب مناسب مدل‌های پایه با روش‌های یادگیری جمعی می‌تواند بهبود چشمگیری در مسائل دسته‌بندی ایجاد کند و به عنوان یک چارچوب قابل اعتماد برای حل مسائل مشابه به کار گرفته شود.

## ۲-۴ - بهبود مدل یادگیری جمعی پشته‌ای وزن‌دار با استفاده از درصد وقوع کلاس‌ها

برای رفع مشکلات موجود در مدل یادگیری جمعی پشته‌ای وزن‌دار، شامل مدل‌های پایه XLNet، RoBERTa، جنگل‌های تصادفی و رأی‌گیری اکثریت، رویکردی جدید مبتنی بر استفاده از درصد وقوع کلاس‌ها به جای پیش‌بینی دقیق کلاس در مدل فرایادگیر درخت تصمیم ارائه شد. این مشکلات شامل انتخاب بیش از یک پاسخ به عنوان بهترین پاسخ یا عدم انتخاب هیچ پاسخی برای برخی سوالات بود. به جای دریافت خروجی‌های صفر و یک از مدل فرایادگیر، خروجی قبل از لایه Softmax که نشان‌دهنده احتمال وقوع هر کلاس است، استفاده شد. این روش به مدل اجازه می‌دهد برای هر پاسخ، درصد تعلق آن به کلاس ۱ (بهترین پاسخ) یا کلاس ۰ (عدم بهترین پاسخ) را محاسبه کند. سپس پاسخ‌ها برای هر سوال بر اساس درصد وقوع مرتب شده و پاسخی با بالاترین احتمال وقوع به عنوان بهترین پاسخ انتخاب می‌شود. این اقدام تضمین می‌کند که برای هر سوال، دقیقاً یک پاسخ به عنوان

قطعی صفر و یک. این رویکرد انعطاف پذیری بیشتری را در تحلیل داده‌ها فراهم می‌کند. مدل از طریق اعتبارسنجی متقاطع مورد آزمایش قرار گرفته و عملکرد آن با دقت بالا بهبود یافته است.

این ساختار ترکیبی و هوشمندانه از قابلیت‌های مدل‌های پیشرفته زبانی، یادگیری ماشین و یادگیری جمعی بهره‌برداری می‌کند تا عملکردی بی‌نظیر در پیش‌بینی بهترین پاسخ‌ها ارائه دهد.

#### ۵- مجموعه داده

برای انجام این پژوهش، از مجموعه داده‌های StackExchange استفاده شده است. این مجموعه داده یکی از منابع جامع و ساختاریافته برای تحلیل داده‌های مرتبط با جوامع پرسش و پاسخ است و شامل اطلاعات متنوعی از سایت‌های مختلف در زیرمجموعه StackExchange مانند Super User، StackOverflow و Android Enthusiasts می‌باشد. داده‌های این مجموعه به صورت فایل‌های XML ارائه شده‌اند و هر فایل حاوی اطلاعاتی درباره پست‌ها، کاربران، نظرات، برچسب‌ها و سایر تعاملات موجود در این پلتفرم‌ها است.

در این پژوهش، از زیرمجموعه android.stackexchange.com استفاده شده است که بر موضوعات مرتبط با اندروید تمرکز دارد. این مجموعه داده شامل پرسش‌ها، پاسخ‌ها و سایر اطلاعات تعاملی کاربران در حوزه توسعه اندروید است و برای تحلیل کیفیت پاسخ‌ها و شناسایی بهترین پاسخ‌ها بسیار مناسب است. داده‌ها به دلیل ساختار دقیق و ویژگی‌های متنوعی که ارائه می‌دهند، امکان استخراج ویژگی‌های متن و غیرمتنی متعددی را برای مدل‌سازی فراهم می‌کنند. انتخاب این مجموعه به دلیل تمرکز آن بر حوزه اندروید و قابلیت تعمیم نتایج به سایر جوامع پرسش و پاسخ انجام شده است.

۱	۰.۹۹	۰.۹۸	۰.۹۹
دقت			۰.۹۹

مدل بهبود یافته با استفاده از درصد وقوع کلاس‌ها توانست مشکلات پیشین را برطرف کرده و دقت و بازخوانی خود را به سطحی بی‌نظیر برساند. با وجود تعداد کمی نمونه که به اشتباه طبقه‌بندی شده‌اند، این روش به دلیل عملکرد قوی در شناسایی بهترین پاسخ‌ها، پایه‌ای مطمئن برای بهبودهای آینده فراهم می‌کند. رویکرد پیشنهادی، نه تنها دقت بالایی ارائه می‌دهد، بلکه روشی تعمیم‌پذیر برای مسائل مشابه در دیگر جوامع پرسش و پاسخ نیز محسوب می‌شود.

بلوک دیاگرام ارائه شده در شکل ۷ ساختار کلی مدل پیشنهادی را برای پیش‌بینی بهترین پاسخ در جوامع پرسش و پاسخ نشان می‌دهد. در این بلوک دیاگرام، داده‌های ورودی ابتدا به دو بخش متنی و غیرمتنی تفکیک می‌شوند. ویژگی‌های متن و غیرمتنی به‌طور جداگانه پیش‌پردازش می‌شوند تا برای مراحل بعدی آماده شوند. این پیش‌پردازش شامل پاک‌سازی داده‌ها، نرمال‌سازی متون و استخراج ویژگی‌های خاص است. سپس داده‌های پیش‌پردازش شده به چهار مدل پایه تغذیه می‌شوند. مدل‌های XLNet و RoBERTa برای استخراج ویژگی‌های متن با کیفیت بالا، جنگل‌های تصادفی برای تحلیل ویژگی‌های غیرمتنی، سایر مدل‌های ترکیبی برای تقویت دقت و افزایش پوشش خطاها مورد استفاده قرار گرفتند. سپس خروجی مدل‌های پایه از طریق یک بلوک ترکیب وزن‌دار به‌طور هوشمندانه وزن‌دهی می‌شود. این فرآیند باعث بهره‌برداری حداکثری از مزایای هر مدل و پوشش ضعف‌های آن‌ها می‌شود. خروجی‌های وزن‌دهی شده مدل‌ها در یک بلوک پشته‌ای ادغام می‌شوند. این مرحله با استفاده از یک متا-یادگیرنده که بر اساس درخت تصمیم‌گیری طراحی شده است، پیش‌بینی نهایی را انجام می‌دهد. برای اطمینان از دقت مدل، پیش‌بینی نهایی بر اساس درصد وقوع خروجی‌ها انجام می‌شود و نه پیش‌بینی

## ۶- جمع‌بندی

در این پژوهش، مدلی کارآمد برای تشخیص بهترین پاسخ‌ها در سیستم‌های پرسش و پاسخ آنلاین با استفاده از روش‌های یادگیری جمعی و تکنیک‌های نوین پردازش زبان طبیعی پیشنهاد شد. پس از پیاده‌سازی و ارزیابی مدل‌های مختلف، نتایج نشان داد که مدل پیشنهادی توانسته است دقت قابل توجهی در تشخیص پاسخ‌های صحیح داشته باشد، اما همچنان چالش‌هایی برای بهبود نتایج وجود دارد. از جمله مهم‌ترین مشکلات، وجود خطاهای FP و FN بود که به دلیل شباهت در درصد وقوع پاسخ‌ها یا انتخاب پاسخ‌هایی با درصد وقوع مشابه به عنوان بهترین پاسخ به وجود می‌آمد.

در فرآیند این تحقیق، مدل‌های مختلفی از جمله XLNet و RoBERTa برای استخراج ویژگی‌های متنی و جنگل‌های تصادفی برای استخراج ویژگی‌های غیرمتنی پیاده‌سازی شدند. در نهایت، با استفاده از روش رأی‌گیری اکثریت و مدل فرا یادگیرنده درخت تصمیم، دقت نهایی به ۹۹.۰۳ درصد رسید. تغییراتی مانند استفاده از خروجی نرم مدل فرا یادگیرنده به جای پیش‌بینی قطعی منجر به بهبود عملکرد در برخی از سوالات شد، به طوری که ۲۹ سوال که در مدل قبلی خطا داشتند، بهبود یافتند.

با این حال، با وجود افزایش دقت، چالش‌هایی مانند انتخاب چندین پاسخ به عنوان بهترین پاسخ یا انتخاب پاسخی که به طور ظاهری بهترین بود اما از نظر کاربر مناسب نبود، همچنان باقی ماند. برای بهبود بیشتر، پیشنهاد می‌شود که خروجی یادگیرنده‌های پایه نیز به صورت درصد وقوع تنظیم شود و این درصدها به مدل فرا یادگیرنده درخت تصمیم داده شود. این اصلاحات می‌تواند به کاهش خطاها کمک کند و دقت مدل را افزایش دهد.

در نهایت، این پژوهش نشان داد که استفاده از مدل‌های پیچیده و یادگیری جمعی در سیستم‌های پرسش و پاسخ می‌تواند به عملکرد بهتری در تشخیص پاسخ‌های صحیح منجر شود، اما نیاز به اصلاحات بیشتر برای رفع چالش‌ها و بهبود دقت دارد. این دستاوردها، زمینه‌ساز تحقیقات آینده در جهت بهبود مدل‌ها و رفع مشکلات موجود خواهند بود.

## ۷- پیشنهادها برای بهبود مدل و تحقیقات آتی

در این پژوهش، مدل یادگیری جمعی پشته‌ای وزن‌دار با استفاده از مدل‌های پایه XLNet، RoBERTa، جنگل‌های تصادفی و رأی‌گیری اکثریت، همراه با فرایادگیر درخت تصمیم مورد بررسی قرار گرفت. اگرچه این روش توانست در بسیاری از موارد عملکرد مناسبی داشته باشد، اما همچنان چالش‌ها و محدودیت‌هایی در شناسایی بهترین پاسخ‌ها وجود داشت. بر این اساس، پیشنهادهایی برای بهبود عملکرد مدل و تحقیقات آتی ارائه می‌شود.

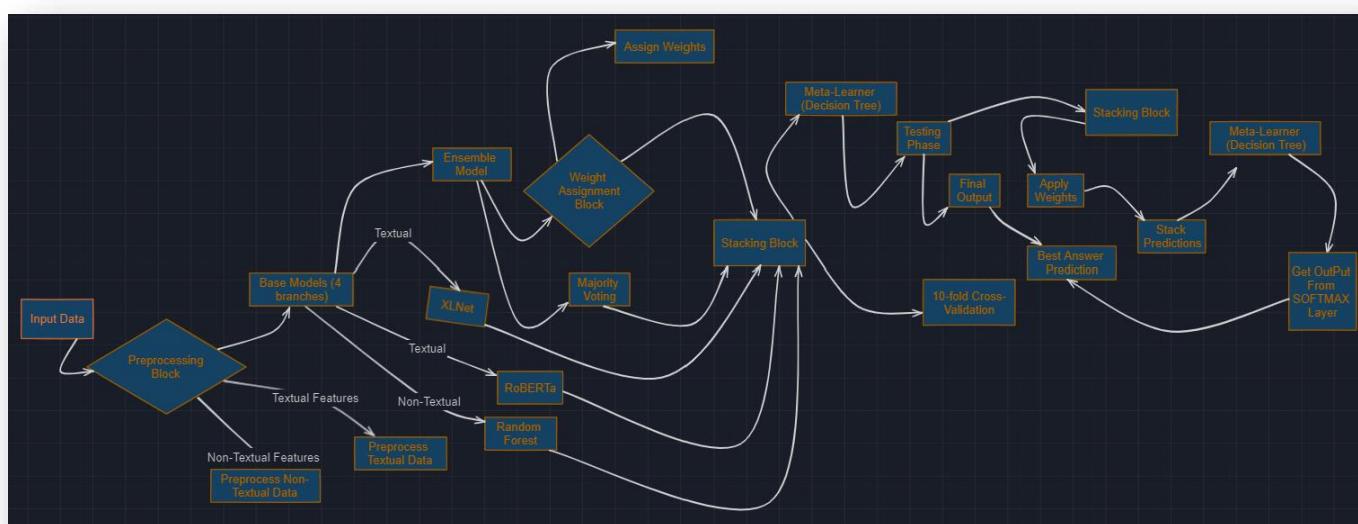
یکی از مهم‌ترین چالش‌های شناسایی شده، استفاده از پیش‌بینی قطعی کلاس در فرایادگیر درخت تصمیم بود. نتایج نشان داد که استفاده از درصد وقوع کلاس‌ها به جای پیش‌بینی قطعی می‌تواند عملکرد مدل را بهبود بخشد و از بروز خطاهای ناشی از همپوشانی احتمالی جلوگیری کند. همچنین، در روش رأی‌گیری اکثریت، وزن‌دهی ثابت مورد استفاده قرار گرفت که ممکن است در شرایط مختلف عملکرد بهینه‌ای نداشته باشد. پیشنهاد می‌شود از رویکرد وزن‌دهی پویا استفاده شود، به طوری که وزن هر مدل پایه بر اساس عملکرد آن در لحظه تنظیم شود. این تغییر می‌تواند به بهبود قابل توجه دقت مدل منجر شود.

تحلیل دقیق‌تر خطاها و استفاده از تکنیک‌های تقویت یادگیری نظیر AdaBoost و Gradient Boosting نیز می‌تواند به کاهش خطاها و یادگیری از اشتباهات گذشته کمک کند. علاوه بر این، گسترش مجموعه داده‌ها با نمونه‌های متنوع‌تر و افزایش تعداد داده‌های آموزشی، به ویژه با افزودن سوالات نامشخص و بدون پاسخ بهینه، می‌تواند مدل را در مواجهه با شرایط واقعی قدرتمندتر سازد. اضافه کردن چنین داده‌هایی، چالش‌های بیشتری برای مدل ایجاد کرده و توانایی آن را در مدیریت شرایط نامطمئن افزایش می‌دهد.

یکی دیگر از رویکردهای پیشنهادی، استفاده از یادگیری انتقالی است. این روش شامل آموزش اولیه مدل بر روی مجموعه داده‌های گسترده‌تر و حاوی اطلاعات پیش‌زمینه‌ای مرتبط است که می‌تواند دانش عمیق‌تری به مدل اضافه کند و توانایی آن را در شناسایی دقیق بهترین پاسخ‌ها افزایش

دهد. برای اطمینان از دقت مدل و کیفیت پاسخ‌های ارائه‌شده، پیشنهاد می‌شود پاسخ‌ها به صورت دستی و بر اساس ۱۳ معیار اصلی ارزیابی شوند که شامل آموزنده بودن، ادب، جامعیت، خوانایی، ارتباط مستقیم با سوال، اختصار، قانع‌کنندگی، جزئیات دقیق، اصالت، واقع‌گرایی، نوآوری، مفید بودن، و کارشناسی‌شدن است. تحلیل نتایج این ارزیابی‌ها با استفاده از ضریب کاپا کوهن می‌تواند میزان تطابق عملکرد مدل با ارزیابی انسانی را به طور دقیق مشخص کند و تصویری شفاف از کیفیت پاسخ‌های تولیدشده ارائه دهد.

در نهایت، برای کاهش تعصبات در سیستم‌های پرسش و پاسخ، الگوریتم‌های کالیبراسیون رأی می‌توانند به کار گرفته شوند. این الگوریتم‌ها با در نظر گرفتن ویژگی‌های رأی‌دهندگان و تعاملات آن‌ها، وزن‌دهی رأی‌ها را بهبود بخشیده و دقت فرآیند انتخاب بهترین پاسخ‌ها را افزایش می‌دهند. اجرای این پیشنهادها می‌تواند به رفع محدودیت‌های کنونی و افزایش کارایی و دقت مدل در شناسایی بهترین پاسخ‌ها منجر شود.



شکل ۷: بلوک دیاگرام کلی مدل



Conference on Computational Linguistics: Technical Papers (pp. ۱۴۲۹-۱۴۴۰).

## مراجع

- [۱۵] Hadfi, R., Moustafa, A., Yoshino, K., & Ito, T. (۲۰۲۲). Best-Answer Prediction in Q&A Sites Using User Information. arXiv preprint arXiv:۲۲۱۲.۰۸۴۷۵.
- [۱۶] Gkotsis, G., Liakata, M., Pedrinaci, C., & Domingue, J. (۲۰۱۵). Leveraging Textual Features for Best Answer Prediction in Community-based Question Answering. arXiv preprint arXiv:۱۵۰۶.۰۲۸۱۶.
- [۱۷] Zhu, Z., Bernhard, D., & Gurevych, I. (۲۰۰۹, November). A multi-dimensional model for assessing the quality of answers in social Q&A sites. In ICIQ (pp. ۲۶۴-۲۶۵).
- [۱۸] Liu, M., Liu, Y., & Yang, Q. (۲۰۱۰). Predicting best answerers for new questions in community question answering. In Web-Age Information Management: ۱۱th International Conference, WAIM ۲۰۱۰, Jiuzhaigou, China, July ۱۵-۱۷, ۲۰۱۰. Proceedings ۱۱ (pp. ۱۲۷-۱۳۸). Springer Berlin Heidelberg.
- [۱۹] Blooma, M. J., Chua, A. Y. K., & Goh, D. H. L. (۲۰۱۰, April). Selection of the best answer in CQA services. In ۲۰۱۰ Seventh International Conference on Information Technology: New Generations (pp. ۵۲۴-۵۳۹). IEEE.
- [۲۰] Burel, G., He, Y., & Alani, H. (۲۰۱۲). Automatic identification of best answers in online enquiry communities. In The Semantic Web: Research and Applications: ۹th Extended Semantic Web Conference, ESWC ۲۰۱۲, Heraklion, Crete, Greece, May ۲۷-۳۱, ۲۰۱۲. Proceedings ۹ (pp. ۵۱۴-۵۲۹). Springer Berlin Heidelberg.
- [۲۱] Gkotsis, G., Stepanyan, K., Pedrinaci, C., Domingue, J., & Liakata, M. (۲۰۱۴, June). It's all in the content: state of the art best answer prediction based on discretisation of shallow linguistic features. In Proceedings of the ۲۰۱۴ ACM conference on Web science (pp. ۲۰۲-۲۱۰).
- [۲۲] Chen, B. C., Dasgupta, A., Wang, X., & Yang, J. (۲۰۱۲, August). Vote calibration in community question-answering systems. In Proceedings of the ۳۰th international ACM SIGIR conference on Research and development in information retrieval (pp. ۷۸۱-۷۹۰).
- [۲۳] Elalfy, D., Gad, W., & Ismail, R. (۲۰۱۸). A hybrid model to predict best answers in question answering communities. Egyptian informatics journal, ۱۹(۱), ۲۱-۳۱.
- [۲۴] Zheng, W., & Li, M. (۲۰۱۷). The best answer prediction by exploiting heterogeneous data on software development Q&A forum. Neurocomputing, ۲۶۹, ۲۱۲-۲۱۹.
- [۲۵] Tondulkar, R., Dubey, M., & Desarkar, M. S. (۲۰۱۸, September). Get me the best: predicting best answerers in community question answering sites. In Proceedings of the ۱۲th ACM Conference on Recommender Systems (pp. ۲۵۱-۲۵۹).
- [۲۶] Chen, D., Fisch, A., Weston, J., & Bordes, A. (۲۰۱۷). Reading wikipedia to answer open-domain questions. arXiv preprint arXiv:۱۷۰۴.۰۰۵۵۱.
- [۲۷] Hu, X. (۲۰۲۳). Enhancing Answer Selection in Community Question Answering with Pre-trained and Large Language Models. arXiv preprint arXiv:۲۳۱۱.۱۷۵۰۲.
- [۲۸] Toba, H., Ming, Z. Y., Adriani, M., & Chua, T. S. (۲۰۱۴). Discovering high quality answers in community question answering archives using a hierarchy of classifiers. Information Sciences, ۲۶۱, ۱۰۱-۱۱۵.
- [۱] Mamykina, L.; Manoim, B.; Mittal, M.; Hripesak, G.; Hartmann, B. Design lessons from the fastest q&a site in the west. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, ۲۰۱۱; pp. ۲۸۵۷-۲۸۶۶.
- [۲] Gkotsis, G.; Liakata, M.; Pedrinaci, C.; Stepanyan, K.; Domingue, J. ACQUA: automated community-based question answering through the discretisation of shallow linguistic features. J. Web Sci. ۲۰۱۵, ۱, ۱-۱۵.
- [۳] Rajpurkar, P.; Zhang, J.; Lopyrev, K.; Liang, P. Squad: ۱۰۰,۰۰۰+ questions for machine comprehension of text. arXiv ۲۰۱۶, arXiv:۱۶۰۶.۰۵۲۵۰.
- [۴] Cui, W.; Xiao, Y.; Wang, H.; Song, Y.; Hwang, S.w.; Wang, W. KBQA: learning question answering over QA corpora and knowledge bases. arXiv ۲۰۱۹, arXiv:۱۹۰۳.۰۲۴۱۹.
- [۵] Zhou, M.; Shi, Z.; Huang, M.; Zhu, X. Knowledge-Aided Open-Domain Question Answering. arXiv ۲۰۲۰, arXiv:۲۰۰۶.۰۵۲۴۴.
- [۶] Correa, D., & Sureka, A. (۲۰۱۴, April). Chaff from the wheat: Characterization and modeling of deleted questions on StackOverflow. In Proceedings of the ۲۲rd international conference on World wide web (pp. ۶۳۱-۶۴۲).
- [۷] Jeon, J., Croft, W. B., Lee, J. H., & Park, S. (۲۰۰۶, August). A framework to predict the quality of answers with non-textual features. In Proceedings of the ۲۹th annual international ACM SIGIR conference on Research and development in information retrieval (pp. ۲۲۸-۲۳۵).
- [۸] Shah, C., & Pomerantz, J. (۲۰۱۰, July). Evaluating and predicting answer quality in community QA. In Proceedings of the ۳۲rd international ACM SIGIR conference on Research and development in information retrieval (pp. ۴۱۱-۴۱۸).
- [۹] Ponzanelli, L., Mocci, A., Bacchelli, A., Lanza, M., & Fullerton, D. (۲۰۱۴, September). Improving low quality StackOverflow post detection. In ۲۰۱۴ IEEE international conference on software maintenance and evolution (pp. ۵۴۱-۵۴۴). IEEE.
- [۱۰] Chen, A., Stanovsky, G., Singh, S., & Gardner, M. (۲۰۱۹, November). Evaluating question answering evaluation. In Proceedings of the ۲nd workshop on machine reading for question answering (pp. ۱۱۹-۱۲۴).
- [۱۱] Hu, Z., Zhang, Z., Yang, H., Chen, Q., & Zuo, D. (۲۰۱۷). A deep learning approach for predicting the quality of online health expert question-answering services. Journal of biomedical informatics, ۷۱, ۲۴۱-۲۵۳.
- [۱۲] Agichtein, E., Castillo, C., Donato, D., Gionis, A., & Mishne, G. (۲۰۰۸, February). Finding high-quality content in social media. In Proceedings of the ۲۰۰۸ international conference on web search and data mining (pp. ۱۸۳-۱۹۴).
- [۱۳] Blooma, M. J., Chua, A. Y., & Goh, D. H. L. (۲۰۰۸, March). A predictive framework for retrieving the best answer. In Proceedings of the ۲۰۰۸ ACM symposium on Applied computing (pp. ۱۱۰۷-۱۱۱۱).
- [۱۴] Suggu, S. P., Goutham, K. N., Chinnakotla, M., & Shrivastava, M. (۲۰۱۶, December). Hand in glove: Deep feature fusion network architectures for answer quality prediction in community question answering. In Proceedings of COLING ۲۰۱۶, the ۲۶th International



- [۲۹] Calefato, F., Lanubile, F., & Novielli, N. (۲۰۱۹). An empirical assessment of best-answer prediction models in technical Q&A sites. *Empirical Software Engineering*, ۲۴, ۸۵۴-۹۰۱.
- [۳۰] Gao, Z., Xia, X., Lo, D., & Grundy, J. (۲۰۲۰). Technical Q&A site answer recommendation via question boosting. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, ۳۰(۱), ۱-۳۴.
- [۳۱] Lou, R., Zhang, K., & Yin, W. (۲۰۲۳). A comprehensive survey on instruction following. *arXiv preprint arXiv:۲۳۰۳.۱۰۴۷۵*.
- [۳۲] Sakai, T., Ishikawa, D., Kando, N., Seki, Y., Kuriyama, K., & Lin, C. Y. (۲۰۱۱, February). Using graded-relevance metrics for evaluating community QA answer selection. In *Proceedings of the fourth ACM international conference on Web search and data mining* (pp. ۱۸۷-۱۹۶).
- [۳۳] Valentin, L. (۲۰۲۲). Answer ranking in Community Question Answering: a deep learning approach. *arXiv preprint arXiv:۲۲۱۲.۰۱۲۱۸*.
- [۳۴] Roy, P. K., Saumya, S., Singh, J. P., Banerjee, S., & Gutub, A. (۲۰۲۳). Analysis of community question-answering issues via machine learning and deep learning: State-of-the-art review. *CAAI Transactions on Intelligence Technology*, ۸(۱), ۹۵-۱۱۷.
- [۳۵] Roy, P. K., Ahmad, Z., Singh, J. P., Alryalat, M. A. A., Rana, N. P., & Dwivedi, Y. K. (۲۰۱۸). Finding and ranking high-quality answers in community question answering sites. *Global Journal of Flexible Systems Management*, ۱۹, ۵۳-۶۸.
- [۳۶] Zoratto, V., Godoy, D., & Aranda, G. N. (۲۰۲۳). A Study on Influential Features for Predicting Best Answers in Community Question-Answering Forums. *Information*, ۱۴(۹), ۴۹۶.
- [۳۷] Mimno, D., Wallach, H., Talley, E., Leenders, M., & McCallum, A. (۲۰۱۱, July). Optimizing semantic coherence in topic models. In *Proceedings of the ۲۰۱۱ conference on empirical methods in natural language processing* (pp. ۲۶۲-۲۷۲).
- [۳۸] Molino, P., Aiello, L. M., & Lops, P. (۲۰۱۶). Social question answering: Textual, user, and network features for best answer prediction. *ACM Transactions on Information Systems (TOIS)*, ۳۵(۱), ۱-۴۰.
- [۳۹] Anbaraki, S., & Jowkar, A. (۲۰۲۱). Evaluating and Predicting the Quality of Answers Factors in the Research Gate's Question and Answer System: a Case Study of the Thematic Domain of Knowledge Management. *Iranian Journal of Information Processing and Management*, ۳۶(۳), ۷۰۹-۷۳۶. doi: ۱۰.۵۲۵۴۷/jipm.۳۶.۳.۷۰۹