

به نام خدا

## تعریف پروژه خبرکاوی

روزانه اخبار بیشماری در تلکس خبرگزاری‌های مختلف قابل مشاهده می‌باشد، حجم این اخبار و گوناگونی آن‌ها به قدری است که از توان پردازش یک فرد یا حتی خود جامعه فراتر می‌باشد. یکی دیگر از معضلات این اخبار ضد و نقیض بودن بعضی از آن‌ها با بعضی دیگر می‌باشد. رایانه در این مسأله می‌تواند به کمک انسان‌ها بیاید. یکی از راه‌های متداول استفاده از سایت‌های خبرخوانی (دیدهبان) می‌باشد به گونه‌ای که خبرها را از سایت‌های مختلف جمع می‌کنند و در یک ساختار مشخص به کاربر نمایش می‌دهند. این سایت‌ها از تکنیک‌هایی از جمله دسته‌بندی و خوشه‌بندی اخبار استفاده می‌کنند تا بتوانند به بهترین نحو ممکن اخبار را به کاربر خود نمایش دهند.

پروژه‌ی خبرکاوی می‌خواهد پای خود را فراتر از این ابزارها بگذارد و بتواند با استفاده از علم داده دانش برخواسته از اخبار را به نمایش بگذارد. این دانش با جمع تمامی اخبار و پردازش بر روی آن‌ها حاصل می‌شود و می‌تواند اطلاعات بسیار مفیدی را در اختیار کاربران قرار دهد. در ادامه به صورت موضوعی خدماتی که این پروژه می‌تواند در اختیار کاربران قرار دهد را می‌آوریم.

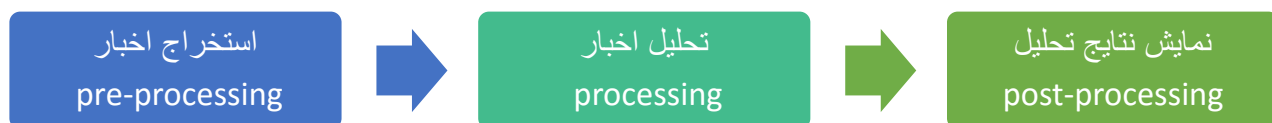
شناسایی انواع مختلف ارتباط بین اخبار یکی از این خدمات می‌باشد. هر دو خبر می‌تواند از طرق مختلف به یکدیگر مرتبط باشند. شباهت دو خبر، هم‌زمانی، دارا بودن تناقض خبری و ... می‌تواند انواعی از ارتباطات بین دو خبر باشند.

شناسایی یک رویداد خبری و جمع تمامی اخبار مرتبط با آن نیز یکی دیگر از خدماتی می‌باشد که در این پروژه در نظر گرفته شده است. رویدادهای خبری را می‌توان به صورت صریح یا ضمنی در متن اخبار مختلف مشاهده کرد. تکنیک‌های هوش مصنوعی می‌تواند به ما کمک کند تا اخباری که دارای اطلاعات جدیدی می‌باشند را به عنوان یک رویداد خبری جدید شناسایی کنند. پس از شناسایی رویدادهای خبری می‌توان تمامی اخبار مرتبط با آن رویداد خبری را جمع کرد و در اختیار کاربر قرار داد تا کاربر بتواند از ابعاد مختلف آن رویداد خبری مطلع شود. علاوه بر این با اضافه کردن بعد زمان به اخبار کاربر می‌تواند سیر اخبار و تغییرات آن رویداد خبری را متوجه شود.

یافتن ارتباط بین اخبار مختلف همچنین این امکان را می‌دهد تا نحوه‌ی انتشار یک خبر در خبرگزاری‌های مختلف نیز به دست آید. این کار کمک می‌کند مرجع یک خبر به خوبی پیدا شود.

## تعریف پروژه

پروژه تحلیل خبر دارای چند قسمت کلی می باشد:



در مرحله‌ی pre-processing تمامی داده‌های خبری استخراج می‌گردند و برای تحلیل آماده می‌گردند. هر خبر به اجزای تشکیل دهنده‌ی آن شامل عنوان، خلاصه، متن خبر، تاریخ خبر و ویژگی‌های دیگری که در هر سایت متفاوت می‌باشد تقسیم می‌شود. این قسمت در نسخه‌ی اول پروژه به صورت آفلاین انجام شده است و خبرهای تعدادی سایت خبری به صورت یکجا<sup>1</sup> دریافت گردیده است ولی در نسخه‌ی آنلاین پروژه به صورت آنلاین با انتشار هر خبر به مجموعه‌ی پایگاه داده‌ی خبری اضافه می‌گردد.

مرحله‌ی processing شامل پیاده سازی و اجرای تمامی الگوریتم‌های تحلیل متن کاوی و خبری می‌باشد و خروجی‌های مورد نیاز برای پروژه را به دست می‌آورد. از جمله کشف ارتباطات بین اخبار مختلف، استخراج کلیدواژه‌ها از آن‌ها، استخراج گراف ارتباط اخبار و سایت‌های خبری و ... که در ادامه‌ی پروژه می‌تواند در بردارنده‌ی قابلیت‌های مختلفی باشد. شناسایی شباهت متن اخبار منتشر شده، شناسایی تناقضات خبری در متن‌هایی که مشابه یکدیگر می‌باشند، استخراج کلیدواژه‌ها و شناسایی رویدادهای خبری از جمله‌ی تحلیل‌هایی هستند که در این مرحله انجام می‌گردند.

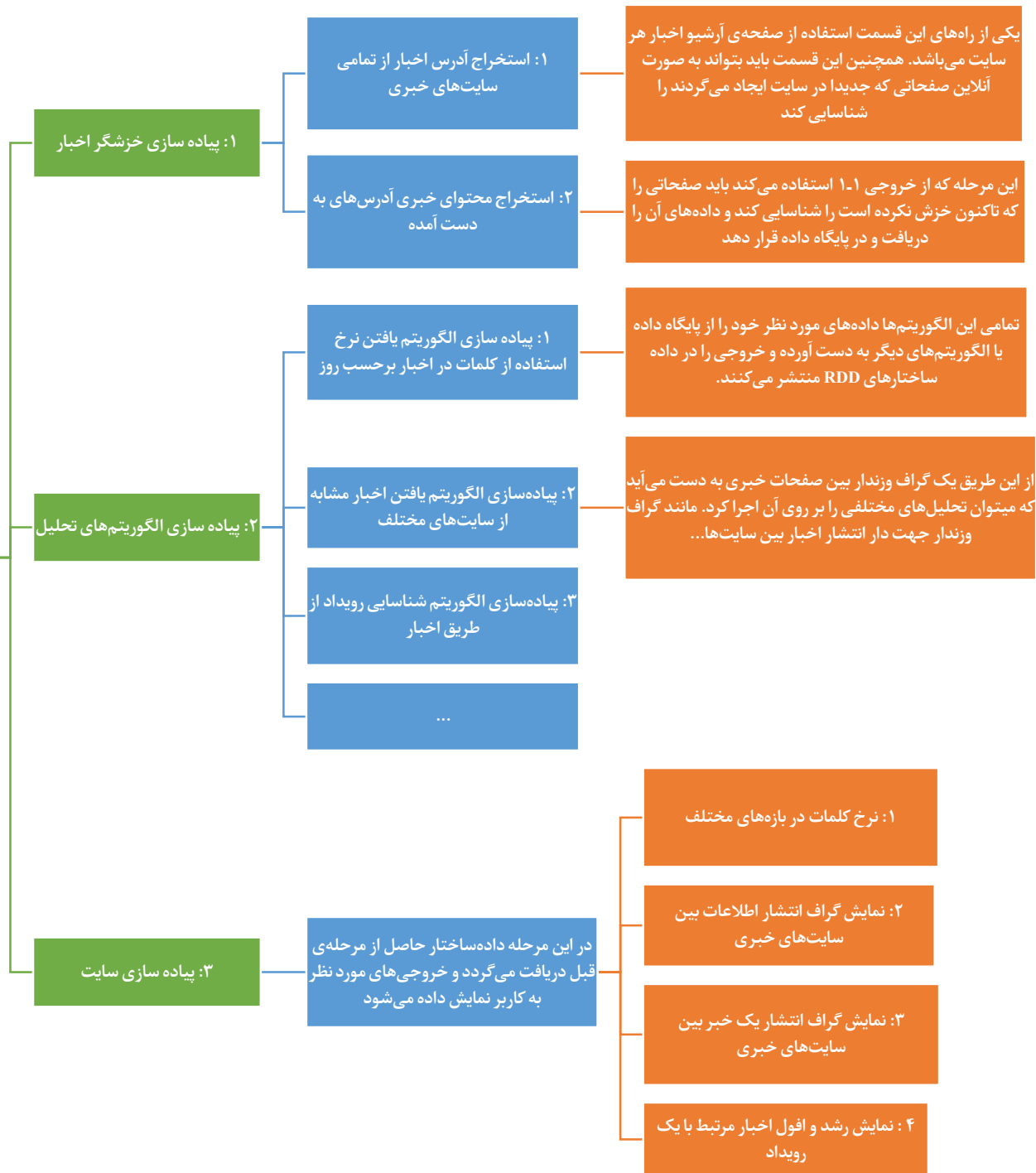
مرحله‌ی post-processing پروژه شامل مورد کاربردهای مختلف می‌باشد که در آن‌ها کاربر می‌تواند به نتایج تحلیل به دست آمده از مرحله‌ی processing دسترسی پیدا کند. از ابزارهای مختلف برای نمایش نتایج استفاده می‌شود. مانند نمودارهای یک، دو یا چند بعدی و گراف‌های مختلف ارتباطی و نمودارهای مختلف آماری.

---

<sup>1</sup> Batch

برای انجام این فازها در ابتدا یک ساختار شکست بین فازهای مختلف پروژه در نظر گرفته شده است:

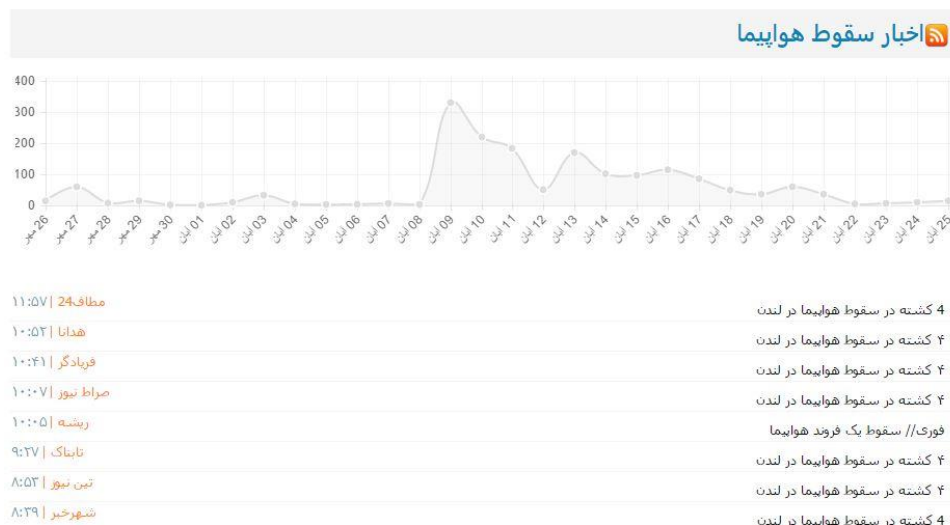
# پروژه تحلیل خبر





نمونه‌ی ایرانی :

• [تی نیوز](#)



این سایت خبری برچسب‌های مهم را شناسایی می‌کند و اخبار مرتبط با آن‌ها نمایش می‌دهد. و همچنین نمودار توزیع برچسب را در دوره‌ی زمانی مشخص می‌کند.

مدل‌های درآمدی برای این گونه پروژه‌ها:

- ارائه سرویس تحلیل خبر خاص منظوره برای سازمان‌ها و نهادهای تصمیم‌گیرنده  
امکاناتی که این ابزار می‌تواند در اختیار سازمان‌ها و نهادهای تصمیم‌گیرنده بگذارد عبارتند از :
  - ✓ ارائه گزارش از موضوعات مختلف در بازه‌ی زمانی مشخص. به این وسیله سازمان می‌تواند در مواردی که مرتبط به فعالیت‌های خود و یا تاثیر گذار بر تصمیمات خود می‌باشد گزارشی دریافت کند که قابل تحلیل و تصمیم‌سازی می‌باشد.
  - ✓ مشاهده‌ی افراد و سازمان‌های مرتبط با یک موضوع و یا یک رویداد خبری. بدین وسیله سازمان می‌تواند این ارتباطات را تحلیل کند و با برقراری ارتباط با آن‌ها دسترسی خود را به موضوع افزایش دهد.
- ارائه سرویس در حوزه‌ی Finance Technology برای مؤسسات مالی و سبدگردان سرمایه به منظور سرمایه‌گذاری کارآمد.

امکاناتی که این ابزار می‌تواند در اختیار این مؤسسات قرار دهد:

- ✓ بررسی وضعیت بازار به وسیله‌ی اخبار مرتبط با بازار (تجمیع اخبار مختلف مرتبط با قیمت‌های خبری)
- ✓ مشاهده‌ی روندهای داغ خبری که می‌توانند در قیمت‌ها تاثیر بگذارند. (پیش‌بینی تغییرات قیمت اقلام و کالاها)
- ✓ شناسایی تاثیر اخبار بر رشد و افول  $EPS^2$  نمادهای بورس.

شناسایی تاثیر اخبار بر رشد و افول  $EPS$  نمادهای بورس.

به عنوان مثال درج خبر ادغام دو شرکت در خبرگزاری‌ها، باید تحلیل شود که این ادغام بر  $EPS$  سهم خواهد افزود یا از آن خواهد کاست که هدف، کسب سود غیر معمول یا جلوگیری از زیان توسط این خبر خواهد بود.

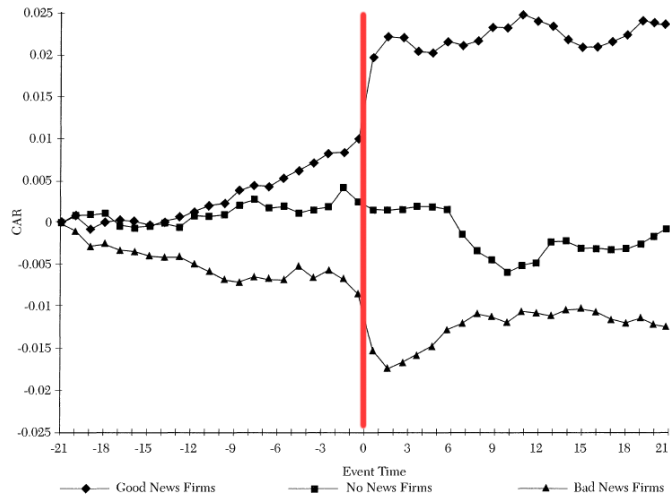
چون مطالعه‌ی رویدادها به بررسی تاثیر یک خبر بر روی جهت و میزان تغییر قیمت یک سهم می‌پردازد، در رشته‌های تحقیقاتی بسیاری می‌تواند کاربرد داشته باشد از جمله در حسابداری، مالی، مدیریت، اقتصاد، بازاریابی، فناوری اطلاعات، حقوق و علوم سیاسی.

برای بررسی تاثیر یک رویداد بر روی یک سهم باید پنجره بر روی گذشته سهم باز کنیم، و تخمین بزنیم که بازده‌ی عادی سهم در روز خبر و چندین روز بعد از آن کجا خواهد بود. بعد از آن، این شیوه به مقایسه بازده‌ی تخمینی و بازده‌ی واقعی برای تعیین میزان بازده‌ی غیر عادی می‌پردازد. در این شیوه، به گذشته‌ی 120 روزه‌ی سهم قبل از خبر نگاه می‌کنیم، و با تحلیل رگرسیون یا استفاده از CAPM یا میانگین بازده تخمینی کوتاه مدت از آینده سهم بدست می‌آوریم و با بازده‌ی واقعی مقایسه می‌کنیم.

برای مشخص کردن اینکه بازده‌ی غیر عادی (اختلاف بین بازده واقعی و تخمینی) غیر صفر است، باید از آزمون‌های آماری مختلف نظیر CAAR- AAR- CAR، -AR- استفاده کرد. یکی از پرکاربردترین تست‌ها، تست t-student است که بازده‌ی غیر عادی را بر خطای جذر میانگین مربعات رگرسیون تقسیم می‌کند.

---

<sup>2</sup> Earnings per share



شکل 1 سه نوع خبر که می‌توانند تاثیر مثبت، منفی یا خنثی در ارزش یک سهم داشته باشند. این دسته‌ها را *Good News*، *Bad News* و *No News* مینامیم