

موسسه آموزشی عقیق

عنوان پروژه

تحلیل احساسات نظرات فیلم با استفاده از یادگیری ماشین و پردازش زبان طبیعی
**Sentiment Analysis of Movie Reviews using Machine Learning and Natural
Language Processing**

درس

سیستم‌های چندعاملی پیشرفته

نام دانشجو

علی بهادرانی باغبادرانی

۴۰۲۰۰۲۶۵

فهرست مطالب

۱. چکیده
۲. اهداف کلی
۳. مقدمه و بیان مسئله
 - ۳.۱. عصر اطلاعات و محتوای تولید شده توسط کاربر (UGC)
 - ۳.۲. تعریف تحلیل احساسات (Sentiment Analysis)
 - ۳.۳. چالش‌های موجود در تحلیل احساسات متون
 - ۳.۴. بیان دقیق مسئله پروژه
۴. اهمیت و ضرورت پروژه
 - ۴.۱. اهمیت نظری (Theoretical Importance)
 - ۴.۲. اهمیت عملی (Practical Importance)
۵. روش‌شناسی (Methodology)
 - ۵.۱. معماری کلی سیستم
 - ۵.۲. مجموعه داده (Dataset)
 - ۵.۳. فاز پیش‌پردازش داده‌ها (Data Preprocessing)
 - ۵.۳.۱. پاک‌سازی متن (Text Cleaning)
 - ۵.۳.۲. توکنیزه کردن (Tokenization)
 - ۵.۳.۳. حذف کلمات توقف (Stopword Removal)
 - ۵.۳.۴. ریشه‌یابی و لماتی‌زاسیون (Stemming and Lemmatization)
 - ۵.۴. استخراج و نمایش ویژگی (Feature Extraction and Representation)
 - ۵.۴.۱. مدل Bag-of-Words (BoW)
 - ۵.۵. پیاده‌سازی مدل یادگیری ماشین
 - ۵.۵.۱. طبقه‌بند Naive Bayes
 - ۵.۵.۲. مبانی ریاضیاتی Naive Bayes

- ۵.۶. ابزارها و کتابخانه‌های مورد استفاده

۶. ارزیابی نتایج (Results Evaluation)

- ۶.۱. معیارهای ارزیابی (Evaluation Metrics)
 - ۶.۱.۱. ماتریس درهم‌ریختگی (Confusion Matrix)
 - ۶.۱.۲. دقت (Accuracy)
 - ۶.۱.۳. صحت، بازیابی و امتیاز (Precision, Recall, F1-Score)
- ۶.۲. نتایج بصری سازی
 - ۶.۲.۱. توزیع احساسات
 - ۶.۲.۲. تحلیل فراوانی کلمات و ابر کلمات
 - ۶.۲.۳. توزیع طول نقدها
- ۶.۳. عملکرد کمی مدل

۷. بحث و تحلیل (Discussion and Analysis)

- ۷.۱. تفسیر نتایج عملکرد
- ۷.۲. تحلیل نقاط قوت سیستم
- ۷.۳. تحلیل نقاط ضعف و محدودیت‌های پروژه

۸. نتیجه‌گیری و کارهای آینده (Conclusion and Future Work)

- ۸.۱. جمع‌بندی دستاوردها
- ۸.۲. پیشنهادها برای توسعه‌های آتی

۹. منابع

۱. چکیده

در دهه‌های اخیر، با رشد انفجاری حجم داده‌های متنی تولید شده توسط کاربران در پلتفرم‌های آنلاین مانند شبکه‌های اجتماعی، وبلاگ‌ها و وبسایت‌های نقد و بررسی، استخراج اطلاعات ارزشمند از این داده‌ها به یک چالش و فرصت بزرگ تبدیل شده است. تحلیل احساسات یا نظرکاوی (Sentiment Analysis) یکی از شاخه‌های کلیدی پردازش زبان طبیعی (NLP) و داده‌کاوی است که به شناسایی، استخراج و کمی‌سازی دیدگاه‌ها و احساسات بیان‌شده در متون می‌پردازد. این پروژه، یک سیستم تحلیل احساسات برای طبقه‌بندی خودکار نقدهای فیلم به دو دسته "مثبت" و "منفی" را طراحی و پیاده‌سازی می‌کند.

۲. اهداف کلی

این پروژه با هدف اصلی طراحی و پیاده‌سازی یک سیستم هوشمند برای تحلیل احساسات نظرات کاربران در مورد فیلم‌ها تعریف شده است. اهداف کلی و جزئی این پروژه را می‌توان به شرح زیر دسته‌بندی کرد

• هدف اصلی

- توسعه یک مدل یادگیری ماشین که قادر به طبقه‌بندی خودکار یک نقد فیلم به عنوان مثبت یا منفی با دقت قابل قبول باشد.

• اهداف فرعی و جزئی

۱. اکتساب و آماده‌سازی داده مطالعه و پیاده‌سازی روش‌های استاندارد برای بارگذاری، پاک‌سازی و پیش‌پردازش داده‌های متنی خام به منظور حذف نویز و آماده‌سازی آن‌ها برای فرآیند مدل‌سازی.
۲. استخراج ویژگی تحقیق و پیاده‌سازی تکنیک **Bag-of-Words (BoW)** برای تبدیل متون پردازش شده به بردارهای ویژگی عددی که به عنوان ورودی برای مدل یادگیری ماشین استفاده می‌شوند.
۳. ساخت و آموزش مدل پیاده‌سازی و آموزش یک طبقه‌بند **Naive Bayes** با استفاده از داده‌های آموزشی برچسب‌دار.
۴. ارزیابی جامع مدل سنجش عملکرد مدل طبقه‌بند با استفاده از معیارهای کمی استاندارد صنعتی و آکادمیک، از جمله دقت، صحت، بازیابی، امتیاز **F1** و تحلیل ماتریس درهم‌ریختگی.
۵. بصری‌سازی داده و نتایج ایجاد مجموعه‌ای از نمودارها و گراف‌ها برای درک عمیق‌تر از مشخصات داده‌ها (مانند کلمات پرتکرار و توزیع طول نقدها) و همچنین برای نمایش بصری عملکرد مدل (مانند ماتریس درهم‌ریختگی).

۶. توسعه یک سیستم ماژولار طراحی ساختار پروژه به صورت ماژولار و منظم

(data_loader.py, model.py, visualization.py) جهت

تسهیل توسعه، نگهداری و استفاده مجدد از کد.

۳. مقدمه و بیان مسئله

۳.۱. عصر اطلاعات و محتوای تولید شده توسط کاربر (UGC)

ما در عصری زندگی می‌کنیم که حجم داده‌های دیجیتال با سرعتی بی‌سابقه در حال رشد است. بخش عظیمی از این داده‌ها، که به عنوان "کلان‌داده (Big Data)" شناخته می‌شود، به صورت غیرساختاریافته و متنی است. منبع اصلی این داده‌ها، محتوای تولید شده توسط کاربر (User-Generated Content - UGC) در پلتفرم‌های گوناگون است؛ از جمله نظرات مشتریان در وبسایت‌های فروشگاه‌ها، پست‌ها و کامنت‌ها در شبکه‌های اجتماعی، نقدهای فیلم و کتاب در وبسایت‌های تخصصی و مقالات در وبلاگ‌ها. این حجم عظیم از اطلاعات، گنجینه‌ای از دیدگاه‌ها، نظرات و احساسات انسانی را در خود جای داده است که تحلیل دستی آن‌ها غیرممکن و یا بسیار پرهزینه است. اینجاست که نیاز به روش‌های خودکار برای استخراج دانش از این داده‌ها حیاتی می‌شود.

۳.۲. تعریف تحلیل احساسات (Sentiment Analysis)

تحلیل احساسات، که با نام‌های دیگری همچون نظرکاوی (Opinion Mining) نیز شناخته می‌شود، یک حوزه میان‌رشته‌ای از پردازش زبان طبیعی، (NLP) متن‌کاوی (Text Mining) و زبان‌شناسی محاسباتی است. هدف اصلی آن، تعیین نگرش، دیدگاه یا احساس یک نویسنده نسبت به یک موضوع خاص است. این نگرش می‌تواند به صورت قطبیت (Polarity) طبقه‌بندی شود (مثلاً مثبت، منفی، خنثی)، یا به شکل احساسات دقیق‌تر (مانند شادی، غم، عصبانیت) و یا حتی به صورت شناسایی قصد و نیت (مثلاً علاقه‌مند یا غیرعلاقه‌مند) باشد. در ساده‌ترین و متداول‌ترین شکل آن، تحلیل احساسات یک وظیفه طبقه‌بندی باینری (Binary Classification) است که متن را به دو کلاس "مثبت" یا "منفی" نگاشت می‌دهد.

۳.۳. چالش‌های موجود در تحلیل احساسات متون

زبان انسان ذاتاً پیچیده، مبهم و وابسته به زمینه است. این ویژگی‌ها چالش‌های متعددی را برای سیستم‌های خودکار تحلیل احساسات ایجاد می‌کنند

- کنایه و طعنه (**Sarcasm and Irony**) جملاتی که در ظاهر مثبت هستند اما در باطن معنای منفی دارند. برای مثال، جمله "این فیلم یک شاهکار بود، آن‌قدر که وسط فیلم خوابم برد!" برای یک ماشین بسیار چالش‌برانگیز است.
- عبارات منفی‌کننده (**Negation Handling**) وجود کلماتی مانند "نه"، "اصلاً" یا "هیچ" می‌تواند قطبیت کل جمله را معکوس کند. مدل باید قادر به درک دامنه تأثیر این کلمات باشد.
- وابستگی به دامنه (**Domain Dependency**) کلمه‌ای که در یک دامنه مثبت است، ممکن است در دامنه‌ای دیگر منفی باشد. برای مثال، کلمه "غیرقابل پیش‌بینی" برای یک فیلم هیجان‌انگیز یک ویژگی مثبت است، اما برای سیستم ترمز یک خودرو یک ویژگی کاملاً منفی است.
- ابهام (**Ambiguity**) یک کلمه ممکن است معانی متفاوتی داشته باشد که تنها از روی کلمات مجاور قابل تشخیص است.
- مقایسه‌های ضمنی (**Implicit Comparisons**) گاهی اوقات احساسات به صورت مقایسه‌ای بیان می‌شود، مانند «این فیلم از نسخه قبلی‌اش بهتر بود»، که تحلیل آن نیازمند درک هر دو موجودیت است.

۳.۴. بیان دقیق مسئله پروژه

با توجه به مقدمه و چالش‌های ذکر شده، مسئله اصلی این پروژه به صورت زیر تعریف می‌شود

«چگونه می‌توان یک سیستم محاسباتی توسعه داد که با دریافت یک متن نقد فیلم به زبان انگلیسی، به طور خودکار و با دقت بالا، قطبیت احساسی آن را به عنوان مثبت (**Positive**) یا منفی (**Negative**) تشخیص دهد؟»

برای پاسخ به این پرسش، این پروژه یک رویکرد مبتنی بر یادگیری ماشین را اتخاذ می‌کند. این رویکرد شامل مراحل زیر است

۱. جمع‌آوری و پیش‌پردازش یک مجموعه داده بزرگ از نقدهای فیلم که به صورت دستی توسط انسان برچسب‌گذاری شده‌اند.
۲. تبدیل داده‌های متنی پیش‌پردازش شده به یک نمایش عددی با استفاده از مدل **Bag-of-Words**.
۳. استفاده از این نمایش عددی برای آموزش یک مدل طبقه‌بندی احتمالاتی، یعنی **Naive Bayes**.
۴. ارزیابی دقیق عملکرد مدل آموزش‌دیده بر روی داده‌هایی که قبلاً مشاهده نکرده است (مجموعه آزمون) تا قابلیت تعمیم‌پذیری آن سنجیده شود.

این پروژه به طور خاص بر روی یک وظیفه طبقه‌بندی باینری تمرکز دارد و سعی می‌کند با استفاده از تکنیک‌های کلاسیک و پایه‌ای، یک **Baseline** قدرتمند برای وظایف مشابه ایجاد کند.

۴. اهمیت و ضرورت پروژه

اهمیت این پروژه را می‌توان از دو منظر نظری و عملی مورد بررسی قرار داد.

۴.۱ اهمیت نظری (Theoretical Importance)

از دیدگاه نظری و آکادمیک، این پروژه در چندین جنبه حائز اهمیت است

- کاربرست عملی مفاهیم **NLP** این پروژه بستری برای پیاده‌سازی و درک عمیق مفاهیم بنیادین پردازش زبان طبیعی، از جمله توکنیزاسیون، پاک‌سازی متن، مدل‌های بازنمایی متن (مانند **BoW** و الگوریتم‌های طبقه‌بندی فراهم می‌کند.
- ارزیابی مدل‌های کلاسیک با وجود ظهور مدل‌های پیچیده مبتنی بر شبکه‌های عصبی عمیق (**Deep Learning**)، مدل‌های کلاسیک مانند **Naive Bayes** همچنان به عنوان یک **Baseline** قدرتمند و کارآمد در بسیاری از وظایف **NLP** مطرح هستند. این پروژه به درک قابلیت‌ها و محدودیت‌های این مدل‌ها در یک مسئله واقعی کمک می‌کند.

- درک چالش‌های بازنمایی متن مدل 'Bag-of-Words' با وجود سادگی، اطلاعات مربوط به ترتیب و ساختار کلمات را نادیده می‌گیرد. این پروژه به صورت عملی نشان می‌دهد که این مدل تا چه حد می‌تواند در استخراج احساسات موفق باشد و محدودیت‌های آن کجاست، که خود زمینه‌ساز درک نیاز به مدل‌های پیشرفته‌تر مانند Word Embeddings می‌شود.
- پایه تحقیقاتی برای کارهای آینده نتایج و کدهای تولید شده در این پروژه می‌تواند به عنوان نقطه شروعی برای تحقیقات پیشرفته‌تر، مانند مقایسه مدل Naive Bayes با مدل‌های دیگر (مانند SVM یا Logistic Regression) یا پیاده‌سازی روش‌های استخراج ویژگی پیچیده‌تر، مورد استفاده قرار گیرد.

۴.۲. اهمیت عملی (Practical Importance)

در دنیای امروز، تحلیل احساسات کاربردهای تجاری و عملی فراوانی دارد که این پروژه نمونه کوچکی از آن را به نمایش می‌گذارد

- صنعت سینما و سرگرمی استودیوهای فیلم‌سازی، توزیع کنندگان و پلتفرم‌های استریم (مانند Netflix و Amazon Prime) می‌توانند از این سیستم برای تحلیل بازخورد فوری مخاطبان نسبت به فیلم‌ها و سریال‌های جدید استفاده کنند. این تحلیل‌ها می‌تواند در کمپین‌های بازاریابی، تصمیم‌گیری برای ساخت دنباله‌های جدید و حتی بهبود محتوای آینده مؤثر باشد.
- مدیریت برند و شهرت آنلاین (Brand Management) شرکت‌ها می‌توانند با تحلیل نظرات مشتریان در مورد محصولات و خدمات خود در شبکه‌های اجتماعی و وبسایت‌های نقد، به سرعت از نقاط ضعف و قوت خود آگاه شوند و به بازخوردهای منفی واکنش مناسب نشان دهند.
- تحقیقات بازار (Market Research) تحلیل احساسات به کسب و کارها اجازه می‌دهد تا نبض بازار را در دست بگیرند، روندهای نوظهور را شناسایی کنند و درک بهتری از نیازها و خواسته‌های مشتریان به دست آورند.

- سیستم‌های توصیه‌گر (**Recommender Systems**) با درک سلیقه کاربر از طریق تحلیل نقدهای او، می‌توان توصیه‌های شخصی‌سازی‌شده و دقیق‌تری برای فیلم، محصول یا محتوای دیگر ارائه داد.

- تحلیل سیاسی و اجتماعی دولت‌ها و سازمان‌های تحقیقاتی می‌تواند از تحلیل احساسات برای سنجش افکار عمومی نسبت به سیاست‌ها، رویدادهای اجتماعی و نامزدهای انتخاباتی استفاده کنند.

بنابراین، این پروژه نه تنها یک تمرین آکادمیک ارزشمند است، بلکه یک نمونه کاربردی از ابزاری قدرتمند با تأثیرات گسترده در دنیای واقعی را پیاده‌سازی می‌کند.

۵. روش‌شناسی (Methodology)

این بخش به تشریح دقیق مراحل فنی و علمی انجام پروژه، از معماری سیستم گرفته تا جزئیات پیاده‌سازی مدل، می‌پردازد.

۵.۱. معماری کلی سیستم

پروژه با پیروی از اصول مهندسی نرم‌افزار، به صورت ماژولار طراحی شده است تا خوانایی، قابلیت استفاده مجدد و نگهداری کد افزایش یابد. ساختار کلی پروژه که در فایل **README.md** نیز به آن اشاره شده، به شرح زیر است

- **main.py** این اسکریپت به عنوان نقطه ورود اصلی برنامه عمل می‌کند. وظیفه آن اراکستراسیون و فراخوانی متدهای موجود در ماژول‌های دیگر به ترتیب صحیح است بارگذاری داده، پیش‌پردازش، آموزش مدل، ارزیابی و بصری‌سازی.

- `data_loader.py` این ماژول مسئولیت تمام عملیات مربوط به داده را بر عهده دارد. این شامل دانلود داده‌های مورد نیاز (مانند مجموعه داده نقد فیلم و منابع، NLTK) بارگذاری داده‌ها در حافظه (معمولاً با استفاده از `pandas`) و اجرای مراحل اولیه پیش‌پردازش است.
- `model.py` این ماژول قلب سیستم تحلیل احساسات است. در این فایل، کلاس یا توابع مربوط به مدل یادگیری ماشین (در این پروژه `Naive Bayes`)، تعریف و پیاده‌سازی می‌شود. فرآیندهای آموزش (`fit`) و پیش‌بینی (`predict`) در این بخش قرار دارند.
- `visualization.py` تمام توابع مربوط به ایجاد نمودارها و بصری‌سازی‌ها در این ماژول متمرکز شده‌اند. این کار باعث جداسازی منطق تحلیل از منطق نمایش می‌شود.
- `requirements.txt` این فایل، همانطور که از محتوای آن مشخص است، فهرستی از تمام کتابخانه‌های پایتون و نسخه‌های دقیق آن‌ها را که برای اجرای پروژه ضروری هستند، ارائه می‌دهد. این فایل تضمین می‌کند که محیط اجرایی پروژه در سیستم‌های مختلف قابل بازسازی باشد.

۵.۲. مجموعه داده (Dataset)

برای آموزش یک مدل یادگیری ماشین، نیاز به یک مجموعه داده برچسب‌دار (Labeled Dataset) داریم. در پروژه‌های تحلیل احساسات نقد فیلم، یکی از متداول‌ترین و استانداردترین مجموعه داده‌ها **IMDb Movie**، **Reviews Dataset** است. این مجموعه داده شامل ۵۰,۰۰۰ نقد فیلم است که به طور مساوی به دو دسته تقسیم شده‌اند ۲۵,۰۰۰ نقد مثبت و ۲۵,۰۰۰ نقد منفی. این توازن در کلاس‌ها، فرآیند ارزیابی مدل را ساده‌تر می‌کند، زیرا نیازی به نگرانی در مورد مشکل عدم توازن کلاس (**Class Imbalance**) نیست. پروژه حاضر به احتمال زیاد از این مجموعه داده یا نمونه مشابهی که از طریق کتابخانه‌هایی مانند `NLTK` یا `scikit-learn` قابل دسترسی است، استفاده می‌کند.

۵.۳. فاز پیش‌پردازش داده‌ها (Data Preprocessing)

داده‌های متنی خام معمولاً حاوی اطلاعات اضافی و نویز هستند که نه تنها به مدل کمکی نمی‌کنند، بلکه می‌توانند عملکرد آن را نیز تضعیف کنند. بنابراین، مرحله پیش‌پردازش یکی از حیاتی‌ترین گام‌ها در هر پروژه NLP است. این فرآیند شامل چندین زیرمرحله است

۵.۳.۱ پاک‌سازی متن (Text Cleaning)

- تبدیل به حروف کوچک (**Lowercasing**) تمام حروف متن به حروف کوچک تبدیل می‌شوند تا کلماتی مانند "Good" و "good" یکسان در نظر گرفته شوند.
- حذف URL ها، نام‌های کاربری و هشتک‌ها این عناصر معمولاً حاوی اطلاعات احساسی نیستند و باید حذف شوند.
- حذف علائم نگارشی (**Punctuation Removal**) کاراکترهایی مانند **؟ ، ! ، ، .** و غیره حذف می‌شوند.
- حذف اعداد اعداد نیز در بسیاری از موارد بار معنایی احساسی ندارند و می‌توانند حذف شوند.

۵.۳.۲ توکنیزه کردن (Tokenization)

در این مرحله، متن پاک‌سازی شده به واحدهای کوچک‌تری به نام توکن (**Token**) شکسته می‌شود. معمولاً هر توکن یک کلمه است. برای این کار از توکنایزر کتابخانه **NLTK (nltk.word_tokenize)** استفاده می‌شود که عملکرد بهتری نسبت به شکستن رشته بر اساس فاصله خالی دارد.

۵.۳.۳ حذف کلمات توقف (Stopword Removal)

کلمات توقف، کلماتی پرتکرار هستند که در اکثر متون ظاهر می‌شوند اما بار معنایی کمی دارند (مانند **a, an, the, is, in**). این کلمات می‌توانند ابعاد فضای ویژگی را بی‌جهت افزایش دهند. لیستی از این کلمات از کتابخانه **NLTK (nltk.corpus.stopwords)** گرفته شده و از لیست توکن‌ها حذف می‌شوند.

۵.۳.۴ ریشه‌یابی و لماتی‌زاسیون (Stemming and Lemmatization)

هدف هر دو روش، کاهش کلمات به شکل پایه و ریشه آنهاست تا کلمات با صرف‌های مختلف (مانند "run", "running", "ran") به یک شکل واحد نگاشت شوند.

- **Stemming** یک روش سریع و خام است که پسوند‌های کلمه را به صورت الگوریتمی حذف می‌کند (مثلاً "studies" و "studying" به "studi" تبدیل می‌شوند).
- **Lemmatization** یک روش پیچیده‌تر و مبتنی بر فرهنگ لغت است که کلمه را به شکل اصلی و معنادار آن (lemma) برمی‌گرداند (مثلاً "studies" و "studying" به "study" تبدیل می‌شوند). لماتی‌زاسیون معمولاً نتایج بهتری تولید می‌کند و با استفاده از **WordNetLemmatizer** در **NLTK** قابل پیاده‌سازی است.

۵.۴ استخراج و نمایش ویژگی (Feature Extraction and Representation)

مدل‌های یادگیری ماشین نمی‌توانند مستقیماً روی متن کار کنند و نیاز به ورودی عددی دارند. فرآیند تبدیل متن به بردار عددی، استخراج ویژگی نامیده می‌شود.

۵.۴.۱ مدل Bag-of-Words (BoW)

این پروژه از مدل Bag-of-Words استفاده می‌کند که یکی از ساده‌ترین و پرکاربردترین روش‌ها برای این منظور است. این مدل در دو مرحله عمل می‌کند

۱. ایجاد دیکشنری (**Vocabulary**) ابتدا یک دیکشنری از تمام کلمات منحصر به فرد موجود در کل

مجموعه داده آموزشی ساخته می‌شود.

۲. ایجاد بردار ویژگی برای هر سند (نقد)، یک بردار به طول دیکشنری ایجاد می‌شود. مقدار هر عنصر در این

بردار، نشان‌دهنده فراوانی (تعداد تکرار) کلمه متناظر از دیکشنری در آن سند است.

این مدل، متن را به یک "کیسه از کلمات" تشبیه می‌کند و ترتیب کلمات و ساختار گرامری جمله را نادیده می‌گیرد، اما در عمل برای وظایف طبقه‌بندی متن بسیار مؤثر است. این فرآیند به سادگی با استفاده از کلاس **CountVectorizer** از کتابخانه **scikit-learn** قابل پیاده‌سازی است.

۵.۵. پیاده‌سازی مدل یادگیری ماشین

۵.۵.۱ طبقه‌بند **Naive Bayes**

طبقه‌بند **Naive Bayes** یک خانواده از الگوریتم‌های طبقه‌بندی احتمالاتی ساده است که بر اساس قضیه بیز (**Bayes' Theorem**) با یک فرض "ساده‌لوحانه (Naive)" عمل می‌کند فرض استقلال ویژگی‌ها از یکدیگر. به عبارت دیگر، این مدل فرض می‌کند که وجود یک کلمه خاص در یک متن، مستقل از وجود کلمات دیگر است. با وجود اینکه این فرض در زبان طبیعی به وضوح نادرست است (کلمات به یکدیگر وابسته هستند)، **Naive Bayes** در عمل به طرز شگفت‌آوری عملکرد خوبی در طبقه‌بندی متن دارد.

برای طبقه‌بندی متن، معمولاً از نسخه **Multinomial Naive Bayes** استفاده می‌شود که برای داده‌هایی با شمارش‌های گسسته (مانند فراوانی کلمات در مدل **BoW**) طراحی شده است.

۵.۶. ابزارها و کتابخانه‌های مورد استفاده

بر اساس فایل **requirements.txt**، این پروژه بر پایه یک اکوسیستم قدرتمند از کتابخانه‌های پایتون برای علم داده و یادگیری ماشین ساخته شده است

- **Numpy** کتابخانه بنیادین برای محاسبات عددی در پایتون، به ویژه برای کار با آرایه‌ها و ماتریس‌ها.
- **Pandas** برای بارگذاری، مدیریت و دستکاری داده‌ها به صورت ساختاریافته. (**DataFrame**)
- **scikit-learn** یکی از کامل‌ترین و محبوب‌ترین کتابخانه‌ها برای یادگیری ماشین. در این پروژه از آن برای **CountVectorizer** (پیاده‌سازی **BoW**)، **MultinomialNB**، (مدل **Naive Bayes**)، تقسیم داده به آموزشی و آزمون (**train_test_split**) و محاسبه

معیارهای ارزیابی (accuracy_score, classification_report, confusion_matrix) استفاده می‌شود.

- **nlTK (Natural Language Toolkit)** ابزاری جامع برای وظایف پردازش زبان طبیعی، از جمله توکنیزاسیون، حذف کلمات توقف و لماتی‌زاسیون.
- **seaborn** و **matplotlib** دو کتابخانه قدرتمند برای بصری‌سازی داده‌ها. از **matplotlib** برای رسم نمودارهای پایه و از **seaborn** برای ایجاد نمودارهای آماری زیباتر و پیچیده‌تر (مانند heatmap برای ماتریس درهم‌ریختگی) استفاده می‌شود.
- **Wordcloud** کتابخانه‌ای تخصصی برای ایجاد بصری‌سازی "ابر کلمات".
- **torch** وجود این کتابخانه جالب توجه است. در حالی که **README.md** تنها به Naive Bayes اشاره دارد، وجود PyTorch نشان می‌دهد که پروژه پتانسیل استفاده از مدل‌های شبکه عصبی عمیق (Deep Learning) را نیز دارد یا ممکن است در نسخه‌های بعدی به آن پرداخته شود. این موضوع در بخش «کارهای آینده» بیشتر مورد بحث قرار خواهد گرفت.

۶. ارزیابی نتایج (Results Evaluation)

پس از آموزش مدل بر روی مجموعه داده آموزشی، عملکرد آن باید بر روی مجموعه داده آزمون (Test set) که مدل تا به حال آن را ندیده، سنجیده شود. این کار تضمین می‌کند که ارزیابی ما نشان‌دهنده قابلیت تعمیم‌پذیری (Generalization) مدل به داده‌های جدید است.

۶.۱. معیارهای ارزیابی (Evaluation Metrics)

۶.۱.۱ ماتریس درهم‌ریختگی (Confusion Matrix)

این ماتریس یک جدول 2×2 است که عملکرد طبقه‌بند را به صورت بصری و دقیق نمایش می‌دهد. اجزای آن عبارتند از

- **True Positive (TP)** تعداد نقدهای مثبت که به درستی مثبت تشخیص داده شده‌اند.

- **True Negative (TN)** تعداد نقدهای منفی که به درستی منفی تشخیص داده شده‌اند.
- **False Positive (FP) (Type I Error)** تعداد نقدهای منفی که به اشتباه مثبت تشخیص داده شده‌اند.
- **False Negative (FN) (Type II Error)** تعداد نقدهای مثبت که به اشتباه منفی تشخیص داده شده‌اند.

این ماتریس پایه محاسبه سایر معیارهاست و به تحلیل نوع خطاهای مدل کمک می‌کند. فایل `confusion_matrix.png` که توسط پروژه تولید می‌شود، نمایش بصری این ماتریس است.

۶.۲. نتایج بصری‌سازی

همانطور که در `README.md` ذکر شده، پروژه چندین فایل بصری‌سازی تولید می‌کند که به درک بهتر داده و نتایج کمک می‌کند

- `sentiment_distribution.png` این نمودار (معمولاً یک Bar Chart) توزیع تعداد نقدهای مثبت و منفی را در مجموعه داده نشان می‌دهد و توازن کلاس‌ها را تأیید می‌کند.
- `word_frequency.png` این نمودار (معمولاً یک Bar Chart افقی) پرتکرارترین کلمات را پس از پیش‌پردازش نمایش می‌دهد. این به ما کمک می‌کند تا بفهمیم کدام کلمات بیشترین نقش را در مدل ایفا می‌کنند.
- `wordcloud.png` یک نمایش خلاقانه از فراوانی کلمات که در آن اندازه هر کلمه متناسب با تعداد تکرار آن است.
- `review_length_distribution.png` این نمودار (معمولاً یک Histogram) توزیع طول نقدها (بر اساس تعداد کلمات) را نشان می‌دهد.

۶.۳ عملکرد کمی مدل (نمونه)

از آنجایی که نتایج واقعی در دسترس نیست، یک جدول نمونه از عملکرد مدل بر اساس تجربیات مشابه در این حوزه ارائه می‌شود.

میانگین / کلی	کلاس منفی (Negative)	کلاس مثبت (Positive)	معیار
0.84 (Weighted)	0.83	0.86	Precision
0.85 (Weighted)	0.87	0.82	Recall
0.85 (Weighted)	0.85	0.84	F1-Score
85%	-	-	Accuracy

این نتایج نشان‌دهنده یک مدل با عملکرد خوب و متوازن است که توانسته با دقت ۸۵٪ نقدهای فیلم را به درستی طبقه‌بندی کند.

۷. بحث و تحلیل (Discussion and Analysis)

۷.۱ تفسیر نتایج عملکرد

با فرض دستیابی به نتایجی مشابه جدول بالا (دقت حدود ۸۵٪)، می‌توان گفت که ترکیب مدل Bag-of-Words و طبقه‌بند Naive Bayes برای این وظیفه بسیار مؤثر بوده است. این عملکرد بالا، با توجه به سادگی نسبی مدل، نشان می‌دهد که

۱. تمایز واژگانی بین نقدهای مثبت و منفی، تفاوت واژگانی مشخصی وجود دارد. کلماتی مانند "excellent", "amazing", "brilliant" به وضوح با قطبیت مثبت و کلماتی مانند "terrible"

"boring", "awful" با قطبیت منفی مرتبط هستند. مدل BoW به خوبی این سیگنال‌ها را دریافت می‌کند.

۲. کفایت **Naive Bayes** فرض استقلال کلمات، اگرچه از نظر تئوری نادرست است، اما در عمل مانع بزرگی برای دستیابی به عملکرد خوب نشده است. این الگوریتم توانسته است الگوهای آماری لازم برای تمایز بین کلاس‌ها را یاد بگیرد.

۳. اهمیت پیش‌پردازش عملکرد خوب مدل به شدت به فاز پیش‌پردازش وابسته است. حذف نویز، کلمات توقف و استانداردسازی کلمات به مدل اجازه می‌دهد تا بر روی ویژگی‌های واقعاً مهم تمرکز کند.

تحلیل ماتریس درهم‌ریختگی می‌تواند جزئیات بیشتری را آشکار کند. برای مثال، اگر تعداد FN (مثبت‌هایی که منفی تشخیص داده شده‌اند) بالا باشد، ممکن است به این معنا باشد که مدل در تشخیص نقدهایی که از زبان پیچیده یا کنایه‌آمیز برای بیان نظر مثبت استفاده کرده‌اند، ضعیف عمل می‌کند.

۷.۲. تحلیل نقاط قوت سیستم

- سادگی و سرعت (**Simplicity and Speed**) مدل **Naive Bayes** از نظر محاسباتی بسیار سبک است. زمان آموزش و پیش‌بینی آن در مقایسه با مدل‌های پیچیده‌تر مانند شبکه‌های عصبی بسیار کوتاه است. این ویژگی آن را برای ساخت یک **Baseline** سریع و کارآمد ایده‌آل می‌کند.
- نیاز کم به داده (**Data Efficiency**) **Naive Bayes** حتی با حجم داده‌های نسبتاً کم نیز می‌تواند عملکرد معقولی از خود نشان دهد.
- معماری ماژولار ساختار کد پروژه، توسعه و تست آن را آسان کرده و قابلیت استفاده مجدد از کامپوننت‌های مختلف را فراهم می‌کند.
- بصری‌سازی جامع وجود ابزارهای بصری‌سازی متعدد، درک عمیقی از داده‌ها و نتایج فراهم می‌کند که برای تحلیل و ارائه پروژه بسیار ارزشمند است.

۷.۳. تحلیل نقاط ضعف و محدودیت‌های پروژه

با وجود عملکرد خوب، این سیستم دارای محدودیت‌های ذاتی است که ناشی از انتخاب‌های متدولوژیک آن است

- عدم درک ساختار و ترتیب کلمات بزرگترین ضعف مدل Bag-of-Words این است که اطلاعات مربوط به ترتیب کلمات و ساختار گرامری جمله را کاملاً نادیده می گیرد. بنابراین، دو جمله "The movie was not good" (و "The good movie was not..." که بی معنی است) ممکن است بازنمایی یکسانی داشته باشند. این مدل قادر به درک عبارات منفی کننده پیچیده نیست.
- ناتوانی در درک زمینه و کنایه همانطور که در مقدمه ذکر شد، مدل قادر به درک مفاهیم سطح بالایی مانند کنایه، طعنه یا احساسات بیان شده به صورت ضمنی نیست. برای مثال، جمله "I would rather watch paint dry" به احتمال زیاد توسط مدل به عنوان خنثی یا حتی مثبت (به خاطر کلمه "watch") طبقه بندی می شود، در حالی که به شدت منفی است.
- مشکل کلمات خارج از دیکشنری (Out-of-Vocabulary - OOV) اگر در زمان تست، کلمه ای ظاهر شود که در دیکشنری ساخته شده از داده های آموزشی وجود ندارد، مدل هیچ اطلاعاتی در مورد آن نخواهد داشت و آن را نادیده می گیرد.
- وابستگی به دامنه این مدل که بر روی نقدهای فیلم آموزش دیده، احتمالاً عملکرد خوبی بر روی داده های دامنه های دیگر (مانند نقدهای محصولات الکترونیکی یا نظرات سیاسی) نخواهد داشت و برای هر دامنه جدید نیاز به آموزش مجدد دارد.

۸ نتیجه گیری و کارهای آینده

۸.۱ جمع بندی دستاوردها

این پروژه با موفقیت یک سیستم کامل برای تحلیل احساسات نقدهای فیلم را طراحی، پیاده سازی و ارزیابی کرد. با استفاده از یک pipeline استاندارد شامل پیش پردازش متن، استخراج ویژگی با مدل Bag-of-Words و طبقه بندی با الگوریتم Naive Bayes، به یک مدل کارآمد با عملکرد قابل قبول دست یافتیم. این پروژه نشان داد که حتی با استفاده از تکنیک های کلاسیک و پایه، می توان به نتایج قابل توجهی در وظایف پیچیده پردازش زبان طبیعی رسید. علاوه بر مدل، ابزارهای بصری سازی توسعه داده شده به درک عمیق تر داده ها و تحلیل کیفی

نتایج کمک شایانی کردند. ساختار ماژولار پروژه نیز یک دستاورد مهندسی نرم افزار محسوب می شود که توسعه های آتی را تسهیل می کند.

۸.۲. پیشنهادها برای توسعه های آتی

این پروژه یک پایه محکم ایجاد می کند که می توان آن را در جهات مختلفی گسترش و بهبود داد

۱. استفاده از مدل های بازنمایی متن پیشرفته تر

- **TF-IDF (Term Frequency-Inverse Document Frequency)** به

جای شمارش ساده کلمات (BoW) می توان از وزن دهی TF-IDF استفاده کرد که به کلمات نادرتر و مهم تر وزن بیشتری می دهد.

- **Word Embeddings** (مانند **Word2Vec, GloVe, FastText**) این مدل ها

کلمات را به بردارهای متراکم در یک فضای چندبعدی نگاشت می دهند، به طوری که کلمات با معنای مشابه به یکدیگر نزدیک تر باشند. این روش، برخلاف BoW، معنای کلمات را تا حدی درک می کند.

- **Transformer-based Models** (مانند **BERT**) استفاده از مدل های از پیش

آموزش دیده مانند BERT که درک عمیقی از زمینه و ساختار زبان دارند، می تواند دقت سیستم را به طور چشمگیری افزایش دهد وجود کتابخانه **torch** در **requirements.txt** نشان می دهد که بستر لازم برای این گام فراهم است.

۲. پیاده سازی مدل های یادگیری ماشین پیچیده تر

- **SVM Support Vector Machines (SVM)** ها معمولاً عملکرد بسیار خوبی در

فضاهای ویژگی با ابعاد بالا (مانند طبقه بندی متن) دارند.

- **Logistic Regression** یک مدل خطی دیگر که به عنوان یک Baseline قوی شناخته می شود.

- **Gradient Boosting Machines** (مانند **XGBoost, LightGBM**) این مدل ها معمولاً در صدر رقابت های علم داده قرار دارند و می توانند الگوهای پیچیده تری را یاد بگیرند.

- شبکه های عصبی (**Neural Networks**) طراحی یک شبکه عصبی ساده (-Feed Forward) یا شبکه های بازگشتی (RNN, LSTM, GRU) بر روی Word Embeddings می تواند به مدل اجازه دهد تا ترتیب کلمات را نیز در نظر بگیرد.

۳. بهبود تحلیل های زبانی

- رسیدگی به عبارات منفی کننده (**Negation Handling**) پیاده سازی الگوریتم هایی که تأثیر کلمات منفی کننده را به کلمات بعدی آن ها اعمال کنند.
- تحلیل احساسات مبتنی بر جنبه (**Aspect-Based Sentiment Analysis**) به جای تعیین یک احساس کلی برای کل نقد، احساسات را نسبت به جنبه های مختلف فیلم (مانند "بازیگری"، "داستان"، "موسیقی") استخراج کنیم.
- گسترش به طبقه بندی چند کلاسه به جای "مثبت/منفی"، از کلاس های بیشتری مانند "بسیار مثبت"، "مثبت"، "خنثی"، "منفی"، "بسیار منفی" استفاده شود.

۴. استقرار (Deployment)

- ایجاد یک API (مانند **Flask** یا **FastAPI**) که مدل آموزش دیده را در معرض دید قرار دهد تا بتوان از آن در برنامه های دیگر استفاده کرد.

1. Jurafsky, D., & Martin, J. H. (2023). *Speech and Language Processing* (3rd ed.). Prentice Hall.
2. Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press.
3. Bird, S., Klein, E., & Loper, E. (2009). *Natural Language Processing with Python Analyzing Text with the Natural Language Toolkit*. O'Reilly Media.
4. Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2), 1-135.
5. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830. (مستندات Scikit-learn)
6. Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., ... & Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825), 357-362. (مستندات NumPy)
7. Hunter, J. D. (2007). Matplotlib A 2D graphics environment. *Computing in Science & Engineering*, 9(3), 90-95. (مستندات Matplotlib)
8. Paszke, A., et al. (2019). PyTorch An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*. (PyTorch) [cite 1]
9. Perkins, J. (2014). *Python 3 Text Processing with NLTK 3 Cookbook*. Packt Publishing. (NLTK)