

لینک کد:

https://github.com/arash-mehrabi-z/data-mining/blob/main/naive_bayes/naive_bayes.ipynb

توضیحات درباره Crawl:

برای درخواست HTTP دادن از کتابخانه Requests پایتون و برای scrape کردن html از BeautifulSoup4 استفاده شده است.
به صورت اتوماتیک به تمام صفحات mstajbakhsh.ir درخواست می‌زند و پست‌های هر صفحه را برمی‌دارد و جداگانه به‌به آن پست‌ها هم درخواست می‌زند.
در پست‌ها نوشته‌هایی که در post-body هستند را به عنوان text بر می‌گردد و کتگوری‌ها را به عنوان class label.

توضیح مدل:

از تابع کتابخانه ای TfidfVectorizer برای درآوردن feature و vectorize کردن اطلاعات استفاده می‌کند و از تابع MultinomialNB برای مدل Naive Bayes.

تحلیل داده‌ها:

من حالت‌های مختلفی را امتحان کردم و نتایج مختلفی را گرفتم که به ترتیب توضیح می‌دهم.

۱. حالت یک:

در داده‌های سایت بعضی text‌ها چندین کلاس مثل (Tutorial, AndroJava و ...) به صورت همزمان دارند و بعضی‌ها تنها یک کلاس.
در این حالت برای تکست‌هایی که چند لیبیل داشت من فقط لیبیل اول را استفاده کردم چون با بررسی چشمی می‌توان نتیجه گرفت که تقریباً همیشه لیبیل اول مرتبط‌ترین به متن پست‌ها بود.
در نتیجه تعداد تکست‌های ما دقیقاً برابر تعداد پست‌های سایت یعنی ۵۵ تا بود.
برای آموزش داده‌ها من از کل داده‌ها استفاده کردم و برای تست کردن ۱۰ پست اخیر وبسایت را انتخاب کردم. در نتیجه دیتای تست را مدل قبلاً دیده بود. نتایج به دست آمده بسیار خوب بود و از ۱۰ تا دیتای ترین ۹ تا را درست تشخیص داده بود.
دیتایی که اشتباه تشخیص داده یک متن مربوط به Home Assistant است که به جای کلاس IoT آن را Anonymity Networks تشخیص داده که احتمالاً به خاطر کم بودن دیتای ترین از کلاس IoT می‌باشد.

۲. حالت دو:

در این حالت همان سیاست استفاده از لیبیل‌های حالت قبلی استفاده شده است.
اما در این حالت دیتای ترین را مدل قبلاً ندیده بود و فقط یکی از ۶ تا دیتای تست را درست تشخیص داده. تستی که درست تشخیص داده مربوط به Hidden Mail Service است که لیبیل Anonymity Networks دارد و دلیل اینکه آن را تشخیص داده احتمالاً به خاطر استفاده از واژه‌هایی است که در دیگر پست‌های مربوط به این کاتگوری هم استفاده شده است مثل Tor, Hidden, server و ...

۳. حالت سه:

در این حالت من همه لیبیل‌ها را استفاده کردم و تکست‌هایی که چندین لیبیل داشتند، به ازای هر لیبیل یک بار دیگر text را وارد دیتاست کردم. پس مثلاً اگر تکست « Einstein was a great scientist » << با دو کاتگوری academic و personal آمده باشد، انگار دو تکست داریم که متن هر کدامشان یکی است اما لیبیل‌هایش متفاوت. طبیعتاً در این حالت چون از لیبیل‌های زیادی استفاده می‌کنیم تعداد داده‌هایمان به جای ۵۵ تا ۱۳۲ تا است و لیبیل‌های پرت زیادی وجود دارد به این معنی که ممکن است از یک لیبیل خاص فقط یک تکست وجود داشته باشد که طبیعتاً عمل‌ترین کردن را خیلی سخت می‌کند. انتظار داریم که در این حالت درصد موفقیت به طور محسوسی پایین‌تر بیاید. برای عمل‌ترین، دوباره مثل حالت یک در این حالت هم مدل قبلاً دیتای تست را دیده.

برای بررسی عمل کرد در این حالت من اینطور بررسی کرده‌ام که از پست‌هایی که چندین لیل دارند حتی اگر یک لیل را هم درست بگویند، آن را به عنوان عمل کرد درست می‌پذیرم. با این فرض می‌بینیم که در این حالت عمل کرد ۱۰۰ درصدی را شاهدیم. علت آن هم این است که طبیعتاً وقتی برای یک تکست چند جواب درست وجود داشته باشد، شانس اینکه مدل بتواند آن را درست تشخیص دهد بالاتر می‌رود.

۴. حالت چهار:

در این حالت همان سیاست استفاده از لیل‌های حالت قبلی استفاده شده است. اما در این حالت دوباره دیتای ترین را مدل قبلاً ندیده بود. می‌بینیم که در این حالت مدل موفق نشده که هیچ کدام از دیتای تست را به درستی تشخیص دهد. به طرز عجیبی می‌بینیم همه دیتای تست را که ۶ پست اخیر وبلاگ هستند را Tutorial پیش‌بینی کرده. اگر به نحوه distribution کاتگوری‌ها نگاه کنیم می‌بینیم که Tutorial با ۱۸ پست بعد از Personal با ۱۹ پست بیشترین تعداد پست‌ها را دارد، اما سایر کاتگوری‌ها مثل AndroJava و Anonymity Networks نیز تعداد خوبی پست دارند پس دیتا نسبتاً خوب توزیع شده و مدل می‌بایست خوب توزیع شده باشد اما عمل کرد مدل در این حالت عجیب است.

نتیجه گیری:

در حالت‌هایی که مدل قبلاً دیتای تست را به عنوان دیتای آموزش دیده بود، عمل کرد بسیار عالی (۹۰٪ و ۱۰۰٪) را داشتیم اما در حالت‌هایی که دیتای تست جدید بود عمل کرد بسیار بد (۱۶٪ و ۰٪) را دیدیم.

دلیل این را هم باید در کمی دیتای آموزش و پخش بودن زیادی کلاس لیل‌ها جستجو کرد. برای بهبود عمل کرد مدل توصیه می‌شود که کاتگوری‌ها را از سایت برداریم بلکه با دیتایی که داریم خودمان تعداد کمتری کلاس لیل (مثلاً ۵ تا) در نظر بگیریم و به صورت supervised به هر تکست یک لیل اختصاص بدهیم. پیش‌بینی می‌کنم که در این حالت عمل کرد مدل در حالت‌هایی که دیتای تست جدید هم باشد بسیار بهتر شود.