



مباحث ویژه در سیستم‌های دیجیتال

مدرس : دکتر ایمان غلامی

[پاییز ۱۴۰۰]

نگارنده: امیرمحمد شعبانی

تمرین سری ۵

مقدمه

سلام. در ابتدا تشکر ویژه‌ای برای آموزش مطالب جالب و کاربردی می‌کنم. امیدوارم درس داده‌های حجیم در سال‌های آینده برای بچه‌های ورودی جدید نیز ارائه شود تا از آن بهره ببرند.

در ادامه روند کلی را توضیح می‌دهم و سپس سعی خود را می‌کنم که هربخش را با جزئیات کامل و تمام چالش‌هایی که وجود داشت توضیح دهم. پس از نصب ملزومات و لود کردن تمام فایل‌ها، سعی کردم با کمی کند و کاو چند نمودار و کمی شهود از داده‌ها بدست بیاورم که بعضی نتایج جالب بودند. سپس سعی کردم بعضی از الگوریتم‌هایی که در طول درس یاد گرفته‌ایم را روی داده‌ها پیاده‌سازی کنم و نتایج را خروجی دهم. اگر چه در پیاده‌سازی بعضی از الگوریتم‌ها ناموفق بودم و یا با وجود پیاده‌سازی درست، در بدست آوردن نتیجه مناسب ناموفق بودم.

همچنین در روند پروژه مشورت‌هایی با خانم ساحل مس فروش در مورد روند پیاده‌سازی و مخصوصاً پیش‌پردازش‌های داده‌های ترافیک و پیدا کردن outliers کرده‌ایم.

نصب

تنها نیازمندی جدی نصب pyspark خواهد بود که در ابتدای برنامه نصب می‌کنیم. البته برای نمایش درست متون فارسی دو کتابخانه دیگر را نصب می‌کنیم و زمان کشیدن نمودارها متن‌ها را با آن‌ها تغییر می‌دهیم.

راه‌اندازی با کولب

بخشی از کدها که برای استخراج داده از درایو و یا وصل شدن به درایو است برای کسانی است که در کولب قصد اجرا آن را دارند. همچنین اگر فایل‌ها در آدرس خاصی قرار دارد کافیهست `base_pass` را ویرایش کنید و زمان لود فایل‌ها به ابتدای آدرس این متغیر را اضافه کنید.

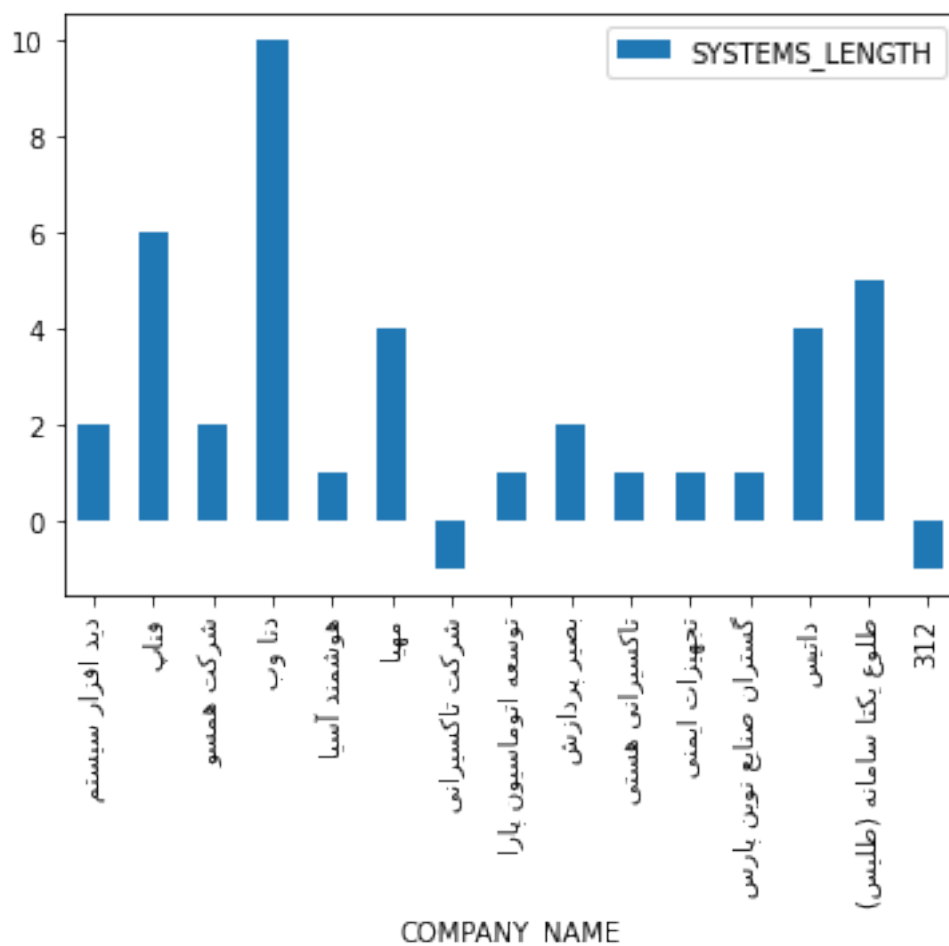
راه اندازی اسپارک و لود فایل ها

اسپارک را روی لوکال می سازیم و سپس schema ترافیک را از روی فایل آن درست می کنیم و سپس داده ترافیک را لود می کنیم. همچنین برای داده های شرکت ها و سیستم ها و راستی آزمایی وضعیت همین کار را اجرا می کنیم. برای بخش هایی که با خط تیره داده ها را جدا کرده است، آن ها را به آرایه تبدیل می کنیم.

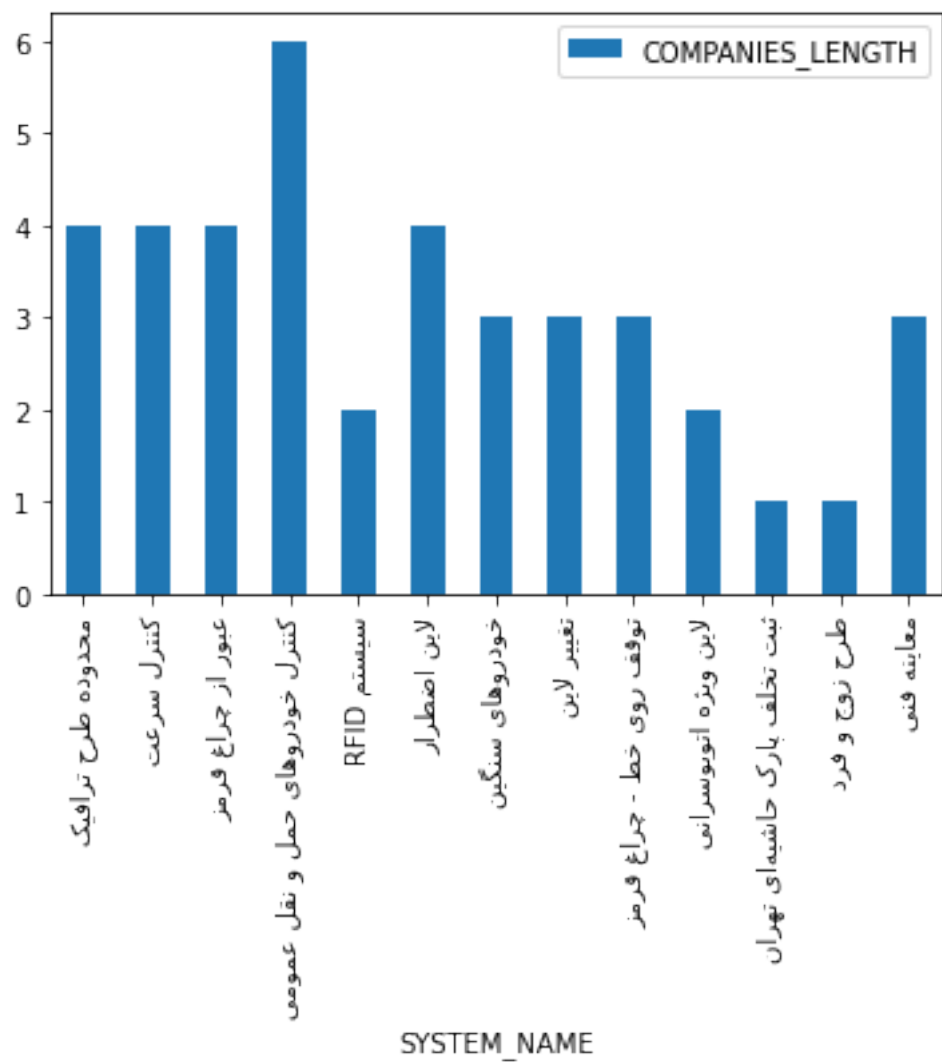
پیش پردازش روی داده ها و بدست آوردن اطلاعات

هر شرکت روی چند سیستم دوربین متفاوت کار می کند؟

هر دوربین برای کار متفاوت ساخته شده است. بعضی از آن ها برای طرح زوج و فرد و بعضی برای چراغ های راهنمایی و رانندگی و ... هستند.

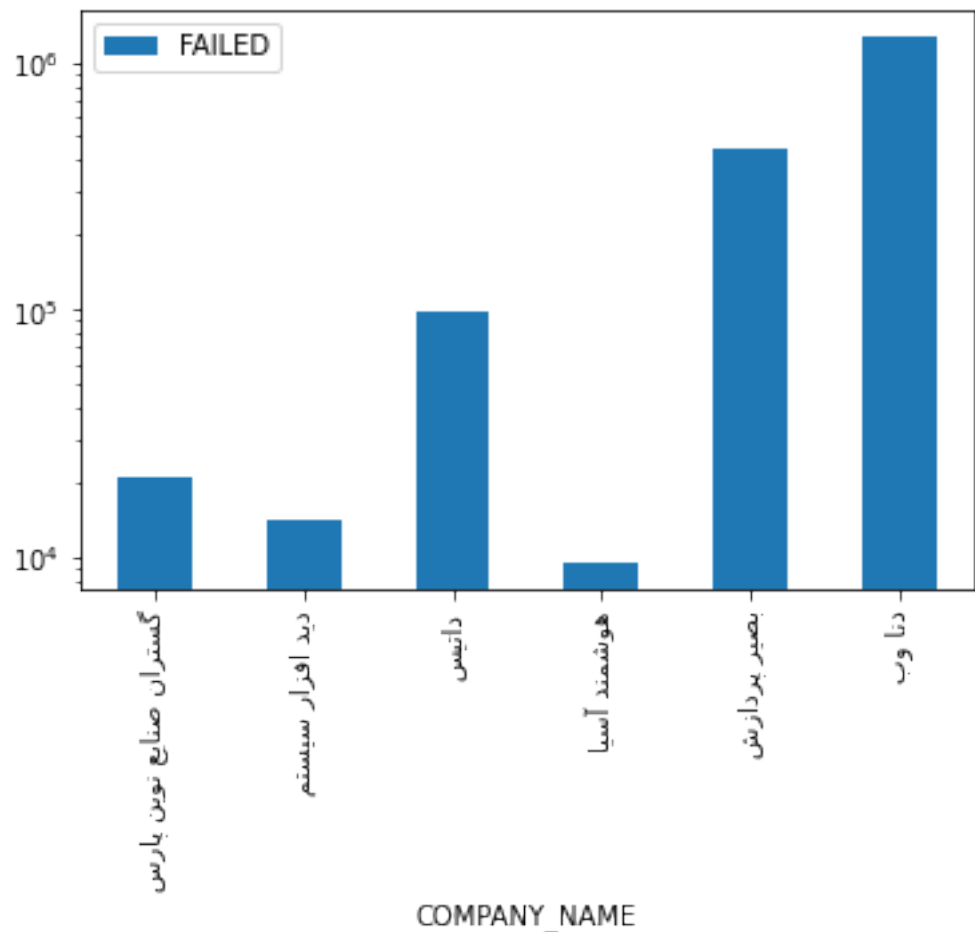


چند شرکت متفاوت روی هرسیستم کار می کند؟



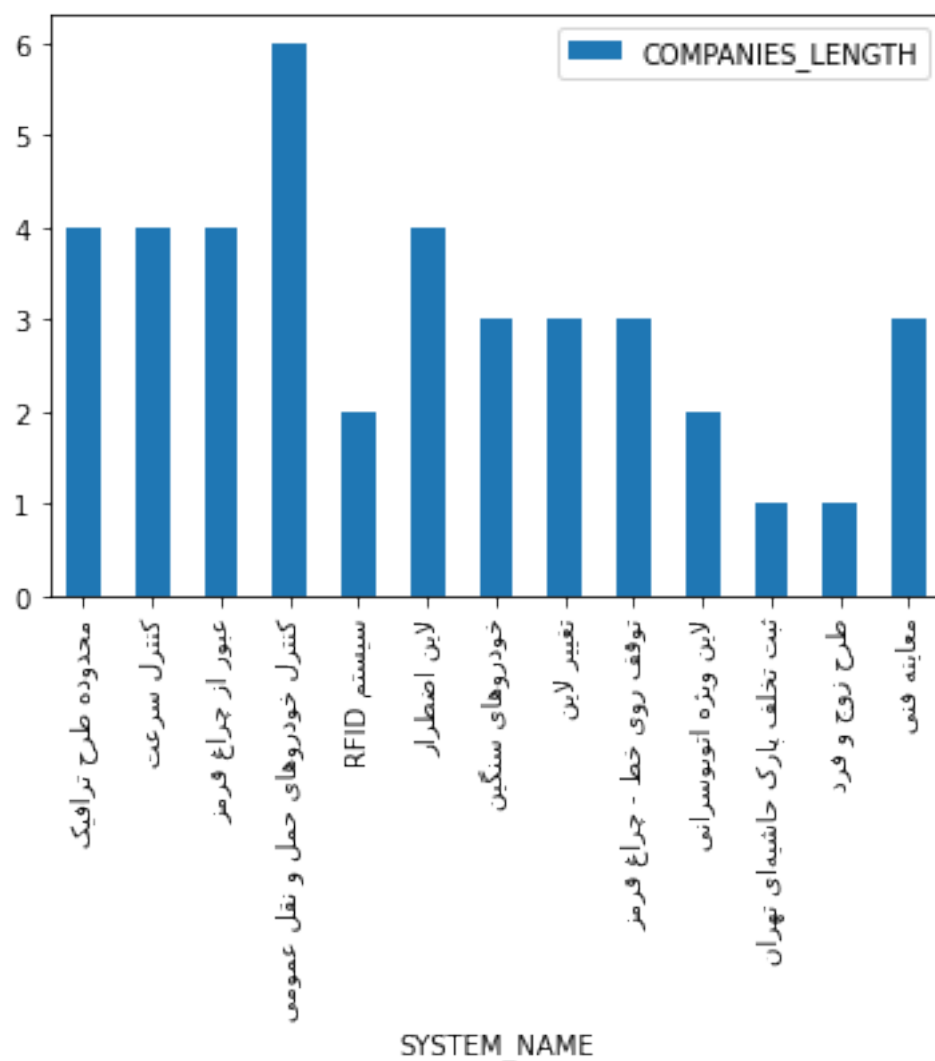
هر شرکت چند بار نتوانست به صورت اتوماتیک پلاک‌ها را تشخیص دهد؟

اگر وضعیت پلاک یک ماشین ۵ یا ۶ باشد آنگاه آن پلاک به صورت اتوماتیک بدست آمده است. بنابراین با فیلتر آن‌ها باقی پلاک‌ها یا درست تشخیص داده نشده‌اند و یا دستی بدست آمده است.



نرخ موفقیت دوربین ها در تشخیص درست یا غلط چقدر است؟

با کمی حساب و کتاب می بینیم هم دوربین هایی وجود دارند که کاملاً درست تشخیص داده اند و هم دوربین هایی وجود دارند که نتوانستند درست تشخیص بدهند. معیار تشخیص را دوربین هایی می گذاریم که نه به صورت اتوماتیک و نه به صورت دستی پلاکشان رهگیری نشده است. سپس تمام آن هایی نرخ تشخیص درستی پلاک آن ها از ۶۰ درصد کمتر باشد را از داده ها حذف می کنیم.



حذف داده‌های پرت

در بخش قبل بخشی از دوربین‌ها را دور انداختیم و حالا نوبت پاک‌سازی بهتر و درست‌تر نسبت به داده‌هایی که داریم است.

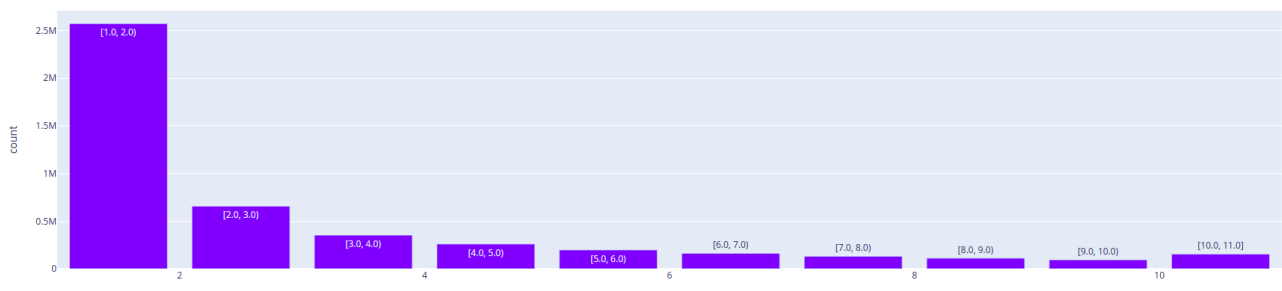
۱۰۰۰ پلاک‌هایی که زیاد دیده شده‌اند و غیر واقعی‌اند

وقتی خلاصه‌ای از داده‌ها را نمایش دادیم میانگین دیده شدن هر پلاک به تقریب ۶ بار بوده است، اما داده‌های کمی وجود داشته‌اند که بیش‌تر از ۱۰ هزار بار دیده شده‌اند. احتمالاً پلاک خودکاری است که دوربین‌ها پس از عدم تشخیص پلاک ثبت می‌کنند. همچنین با ایمیل از شما نیز پرسیدیم و تایید کردید. بنابراین تصمیم گرفتم چنین پلاک‌هایی را حذف کنم و تصمیم بر این شد پلاک‌هایی که کمتر از ۱۰۰ بار دیده شده‌اند را نگه دارم. این تصمیم را با دیدن داده‌ها و با احتساب اینکه تنها برای ۸ روز است انتخاب کردم. همچنین راه مناسب دیگر این است که چارک اول و سوم را محاسبه کنم و یک بازه مناسب با استفاده از آن‌ها بسازم و تنها پلاک‌هایی که به همان تعداد دیده شده‌اند را بگیرم.

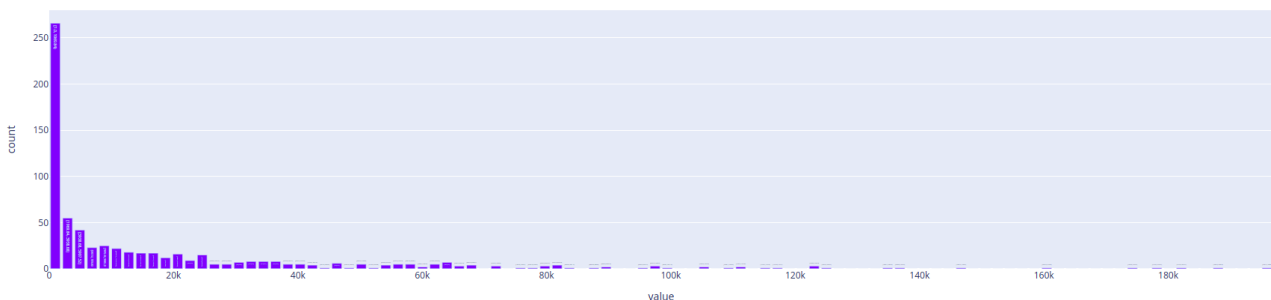
بدست آوردن پلاک‌ها و دوربین‌های پرتدد

کافیست یک بار براساس پلاک و بار دیگر براساس دوربین تعداد بار دیده‌شدن را بشماریم و سپس با کشیدن هیستوگرام آن پرتدد ها را تشخیص دهیم و سپس فیلتر کنیم. پلاک‌هایی که بیشتر از ۷ بار دیده‌اند پرتدد محسوب می‌شوند و دوربین‌هایی که بیشتر از ۸۰ هزار ریکورد دارند را پرتدد می‌نامیم.

هیستوگرام شمارش دوربین‌ها



هیستوگرام شمارش پلاک‌ها



Collabarating Filtering

در بیشتر الگوریتم‌ها هدف بهینه‌بودن کد و تا جای ممکن داشتن زمان مناسب بوده‌است. برای اینکار سعی شده از توابع آماده خود pyspark استفاده کنم چون درغیراینصورت بسیار کند خواهد بود و الگوریتم منطقی نخواهد بود؛ زیرا مسئله زمان پردازش در داده‌های حجیم زیر سوال می‌رود.

با کمی جست‌وجو به الگوریتم ALS در اسپارک می‌رسیم که با استفاده از کم‌ترین مربعات و یادگیری ماشین سعی بر پیاده‌سازی سریع الگوریتم CF را دارد. در این الگوریتم پلاک‌ها حکم کاربرها و دوربین‌ها حکم آیتم‌ها را برای ما دارند و تعداد باری که هم را دیده‌اند حکم امتیازی که یک پلاک به یک آیتم می‌دهد را دارد. داده‌ها را به دو بخش تمرین و تست تقسیم می‌کنیم و سپس سعی می‌کنیم الگوریتم را اجرا کنیم. مقدار خطای آن ۰.۹۲ بوده‌است. در ادامه چندتا از پیش‌بینی‌ها برای داده‌های تست را نشان داده‌ام و سپس برای هر دوربین چند پلاک که حتما می‌بیندش با پیش‌بینی تعداد باری که می‌بیند و برای هرپلاک چند دوربین با پیش‌بینی تعداد باری که می‌بیند را گذاشته‌ام.

تقسیم داده به بازه‌های ۶ ساعت

برای گرفتن پیشنهاد با استفاده از بازه‌های ساعتی و پلاک‌ها یا بازه‌های ساعتی و دوربین‌ها نیازمند تبدیل داده‌ها بودیم. ابتدا روز زمان ثبت را در ستونی ذخیره می‌کردم و سپس ساعت‌ها را با تقسیم بر ۶ کردن به اندیس‌های به خصوصی مپ کردم و با ساختن رشته جدید با چسباندن این دو توانستم رکوردها را با استفاده از پلاک و دوربین و ساعت و روز به خصوص گروه‌بندی کنم و سپس شمارش را انجام دهم.

از آنجایی که الگوریتم‌هایی که از آن‌ها استفاده می‌کنیم نیازمند یک عدد و نه رشته برای هرکدام بود باید آن‌ها به ایندکسی مپ می‌کردم. سپس یک بار توسط پلاک و زمان و بار دیگر توسط دوربین و زمان گروه‌بندی کرده‌ام و شمارش کردم.

Colaborating Filtering Per 6 Hour

Device

همان الگوریتمی که در بخش‌های قبل توضیح داده‌ام را استفاده می‌کنیم. با این تفاوت که کاربرها همان دوربین‌ها هستند و آیتم‌ها بازه‌های ساعتی ۶ ساعته هستند.

Plate

همان الگوریتمی که در بخش‌های قبل توضیح داده‌ام را استفاده می‌کنیم. با این تفاوت که کاربرها همان پلاک‌ها هستند و آیتم‌ها بازه‌های ساعتی ۶ ساعته هستند.

Utility Matrix & SVD

بخشی که به شدت چالش برانگیز بود و بسیار نیازمند فکر زیاد بود، ساختن ماتریسی Utility از رکوردها بود. خوش‌بختانه با جست‌وجو برای مطالب دیگر به طور اتفاقی روشی برای ساخت این ماتریس در زمان معقول پیدا کردم. ماتریسی که ساخته می‌شود باید هرسطر برای یک پلاک باشد که هراندیس آن نشان‌دهنده تعداد باری است که آن دوربین به خصوص را دیده‌است. اما برای اینکار یک سری از زوج‌های دوربین‌ها و پلاک‌ها را نداشتیم و باید کراس می‌کردیم و اینگونه داده اضافی تولید می‌شد. همچنین برای ساخت بردار باید ترتیب‌ها در همه سطرها درست باشد و به عبارتی باید تابع سورت را اجرا می‌کردیم که باعث زیاد شدن زمان اجرا می‌شد. بنابراین از ایندکس‌گذاری روی تمام دوربین‌ها و همچنین طریقه ساخت بردار اسپارس کمک گرفتم و برای هررکورد ایندکس دوربین را به تعداد باری که دیده‌شده مپ کردم و سپس با استفاده از پلاک گروه‌بندی کردم و توسط این مپ‌ها بردار اسپارس را ساختم. البته راه آسان‌تر که بعداً به ذهنم رسید این بود که یک لیست از ایندکس‌ها و یک لیست از شمارش‌ها بسازم و آن‌دو را به بردار اسپارس بدهم تا برایم بسازد. حال که بردارها را داریم کافیهست که آن را به RowMatrix بدهیم تا بتوانیم از طریق آن svd آن را محاسبه کنیم.

پس از محاسبه ماتریس را بر بردارهای V و U ضرب کردم اما هیچ نتیجه خاصی پیدا نکردم. همچنین همین کار را با تعداد فاکتورهای متفاوت کردم و نتوانستم به داده مناسب خوبی برسم. فاکتورها ۵ و ۱۰ و ۱۵ و ۲۰ را امتحان کردم.

Clustering

با استفاده از برداری که بدست آورده ایم کلاستر کردن داده ها نیز ساده خواهد بود اما چون الگوریتمی که از آن استفاده می کنیم برای پکیج یادگیری ماشین است باید بردار همان را استفاده کنیم برای همین بردارها به بردار یادگیری ماشین تبدیل کرده ام.

البته شایان ذکر است که در این بخش آنقدر کار نکرده ام و چیز زیادی برای ارائه ندارم. تنها از الگوریتم دریکله که در کلاس ارائه شده استفاده کرده ام و سعی کردم تاپیک ها را برای هریخش بدست آوردم اما از آنجایی که نتوانستم شهود خوبی در بخش svd روی فاکتورها پیدا کنم، در این قسمت نیز شکست خوردم و تحقیق بیشتر را بعد از نتیجه گرفتن در svd گذاشته ام که متأسفانه به آن هم نرسیدم.

مسائل مربوط به گراف

پس از آن سعی کردم الگوریتم های گرافی را تست کنم و سعی کنم اطلاعات باارزشی بدست آوردم اما متأسفانه با سرچ متوجه شدم که پکیج graphx در اسپارک جدید است و تنها اسکالا آن را ساپورت می کند و برای همین باید با اسکالا پیاده سازی می شد. اما متأسفانه به دلیل مشکلات شخصی مانند اوایل نتوانستم روی پروژه کار کنم و روی اسکالا ادامه ندادم.

البته به این هم فکر کردم که خودم آن را روی پایتون خام پیاده سازی کنم اما هم حلقه در پایتون کند است و هم داده ها بسیار زیاد است و این کار منطقی نبود. چون پیاده سازی آن ها آنقدر سخت نبود و همانطور که اجازه داده بودید به شخصه ظرف چند دقیقه با کد برای یکی از سوال ها کد آن را زدم ولی هدف پیاده سازی برای داده های حجیم بود که اینکار در پایتون و بدون کتابخانه درست منطقی نبود.

Streaming

در این بخش الگوریتمی روی داده ها برای ارائه ندارم و صرفاً کدهایی که با آن ها روی مسلط شدن بر بخش استریم اسپارک کار می کردم قرار دارد. دلیل اینکه به مشکل خوردم این بود که وقتی یک کد را اسپارک اجرا می کند تمام آن بخش با تمام حافظه را بین تمام ورکرها کپی می کند و موازی جلو می برند؛ بنابراین در زمان استریم اگر یک دیکشنری که برای هر ای دی نگه داشته که چندبار دیده را ذخیره کرده باشیم، در زمان آمدن داده های جدید دیکشنری ریست می شود و به مقدار اولیه قبل از اجرای کد می رود و به همین دلیل شمارش اشتباهی صورت می گیرد و به همین منظور برای هر داده نمی توانیم نگه داریم که چندمین بار است که دیده می شود و به همین دلیل نمی توانیم مومنت ها را بشماریم.