



دانشکده مهندسی برق

مباحث ویژه در سیستم‌های دیجیتال

نیم‌سال اول ۱۴۰۰-۱۴۰۱

مدرس: دکتر ایمان غلام‌پور

توضیحات پروژه

شماره دانشجویی: ۹۷۱۰۱۰۲۶

نام و نام‌خانوادگی: امین کشیری

۱ مقدمه

در این پروژه، با استفاده از ایده‌ها و روش‌های مختلفی که در طول ترم آموختیم، اطلاعاتی را از داده‌ها بیرون کشیدیم. بخش‌های مختلف این پروژه را در فایل‌های jupyter جداگانه قرار داده‌ام، و توضیحات هر بخش را نیز جداگانه در این فایل نوشته‌ام.

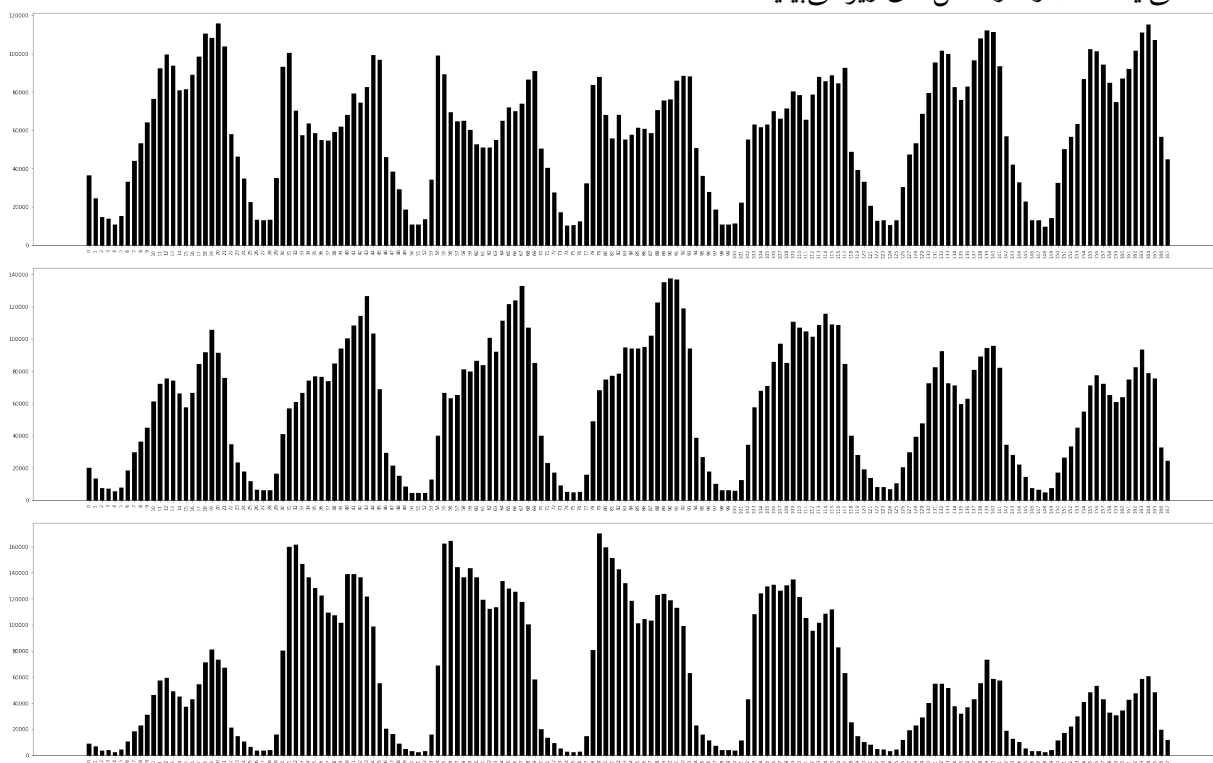
۱.۱ توضیحات کلی

۱. در ابتدای تمامی کدها تنظیمات اولیه اسپارک را انجام دادم، و سپس فایل csv داده شده را load کردم.
۲. در بعضی از بخش‌ها، روز ۸م را از داده‌ها حذف کردم. دلیل این کار این بود از تمامی روزها به اندازه‌ی متناسب با هم داده داشته باشیم. در غیر این صورت تعداد داده‌ها از روز سه شنبه دو برابر باقی روزها می‌شد. کارهای دیگری نیز می‌توانست انجام بگیرد. مثلاً می‌شود داده‌های روز سه شنبه را میانگین بگیریم (یعنی در تمام قسمت‌هایی که تعداد متغیری را شمرده‌ایم، برای روز سه شنبه این تعداد را تقسیم بر ۲ کنیم). در بعضی از قسمت‌ها اما این زیادتر بودن داده‌های روز سه شنبه مشکلی ایجاد نمی‌کرد. اما دقت کنید که در بعضی از قسمت‌های دیگر می‌تواند تحلیل ما را دچار انحراف کند (مثلاً ممکن است به اشتباه نتیجه بگیریم که روز سه شنبه روز پر تردد تری است، یا دوربین‌هایی که در روز سه شنبه دیده می‌شوند را به اشتباه مهم تر در نظر بگیریم).
۳. توضیحات کد و روند اجرا را در فایل‌های jupyter نوشته‌ام. سعی کرده‌ام که توضیحات منطق پشت کدها را در این مستند بنویسم (و نه در خود کدها). بنابراین توضیحات تکنیکال خود کد در اینجا کمتر نوشته شده است.

۳ Clustering

۱.۳ توضیحات کلی

با استفاده از خوشه سازی، سعی کردم دوربین ها را به دسته های مختلفی تقسیم کنم، و برای هر دسته مفهومی بیابم. نماینده هر دوربین در این روش، یک بردار با اندازه 24×7 است، که در هر خانه آن تعداد تردد در آن ساعت از روزهفته قرار گرفته است. ۲۴ ساعت اول برابر با یکشنبه است، ۲۴ ساعت بعدی برای دوشنبه و الی آخر. پس بردار متناظر هر دوربین، تعداد تردها در هر ساعت از یک هفته را برای آن دوربین مشخص می کند. برای خوشه سازی از الگوریتم LDA یا Latent Dirichlet Allocation استفاده شده است، که همانطوری که در درس دیدیم استفاده اولیه آن پیدا کردن توزیع topic های مختلف و کلمات آن ها برای هر مقاله است. با استفاده از این الگوریتم برای داده های تردد ماشین ها نیز می توانیم دقیقاً به چنین توزیعی برسیم. یکی از متغیرهای بسیار مهم در این بخش، `cluster_center` است، که تعداد کلاسترهای نهایی را مشخص می کند. با تغییر این متغیر این متغیر می توانیم تعابیر متفاوتی از داده داشته باشیم. اما یکی از واضح ترین نتیجه ها برای `cluster_center = 3` به دست می آید که آن را در شکل های زیر می بینید:



۲.۳ تحلیل نتایج

مهم ترین نکته ی این سه تصویر، روند تغییر تردها در هر روز است. در دسته ی اول، تردد در ساعات اولیه روز افزایش می یابد، در هنگام ظهر کاهش پیدا می کند، سپس دوباره در شب افزایش پیدا می کند. در دسته ی دوم، تردد در صبح کم است، اما کم کم افزایش می یابد و در شب به اوج خود می رسد. دسته ی سوم روندی دقیقاً عکس دسته ی دوم دارد، و بیشترین تردد را در صبح دارند و سپس کاهش می یابد.

سه دسته ی بالا را می توانیم به این صورت تفسیر کنیم. دسته ی اول نقاط پر تردد شهر هستند، که هم در روز و هم در شب تردد بالایی دارند. این نقاط احتمالاً مکان هایی وسط شهر هستند که تمام طول روز تردد دارند (البته طبیعتاً تردد در ظهر کاهش می یابد). دسته ی دوم، احتمالاً مکان های دیدنی و تفریحی و یا بازارهای شبانه هستند که در طول روز تردد زیادی ندارند (به دلیل این که مردم مشغول کار و مدرسه و ... هستند). دسته ی سوم نیز احتمالاً مکان هایی هستند که در طول روز تردد بالایی دارند، مانند مکان های اداری، مسیر مدارس و ادارات، یا دوربین های نزدیک به مثلاً نانوائی ها.

یک نکته ی بسیار جالب دیگری که در این تصاویر دیده می شود، این است که تردد در روزهای جمعه، شنبه و یکشنبه به طرز جالبی پایین است. با چک کردن این ۳ روز روی تقویم، فهمیدم که این روزها تعطیل رسمی بوده اند (قیام ۱۵ خرداد و شهادت امام جعفر صادق (ع)). بسیار جالب است که این کاهش تردها، فقط در دسته ی سوم رخ داده است، که

دقیقا با شهود ما همخوانی دارد، که دسته‌ی سوم مکان‌هایی مانند مدارس و ادارات هستند. همچنین، الگوی کاهشی تردد در این سه روز از بین رفته است که باز هم مطابق با الگوی پیدا شده است.

