

(۱) فهم مسئله:

۱. انگیزه اصلی چنین پروژه‌های چیست؟ اهداف کسب و کار انگیزه اصلی پروژه است. اهدافی مانند این که ویژگی‌های مشترک مشتریانی که اخیراً از دست داده ایم و از خدمات و محصولات شرکت‌های رقیب استفاده می‌کنند چیست؟ و یا هریک از مشتریان شرکت دارای چه ارزشی برای ما هستند.
۲. خروجی چنین پروژه‌های برای چه مواردی ممکن است کاربرد داشته باشد؟۱- اینکه یک مقاله خبری به چه دسته‌ای تعلق دارد۲- پیش بینی حجم فروش برای تاریخ معین آینده۳- رفتارهای مشتریان را میتوان مشخص کرد
۳. چه کسانی ممکن است به نتایج این پروژه علاقمند باشند؟ چرا؟ صاحبان کسب و کار(پیش بینی نرخ فروش...)

(۲) درک داده:

۱. داده‌ها از کجا بدست آمده اند و چگونه جمع آوری شده اند؟ براساس یک فایل اکسل که مقادیر و مشخصات قیمت خانه‌ها در یک منطقه را بررسی میکند
۲. هر یک از متغیرها چه چیزی را اندازه گیری میکنند؟ متغیرهای شماره ۱ تا ۵ (مستقل) هستند و متغیر شماره ۶ (price) متغیر وابسته است. براساس محور X مقادیر محور Y مشخص میشوند یعنی محور X متغیر مستقل است و متغیر Y وابسته است. و متغیر address را هم حذف کردیم چون نوع آن object بود. در مدل زدن باید متغیرهای وابسته و مستقل مشخص شوند با دستور iloc مشخص میکنیم کدام وابسته و کدام مستقل باشد.
- ۱) 'Avg. Area Income': Avg. Income of residents of the city house is located in. درآمد متوسط ساکنان منطقه.
- ۲) 'Avg. Area House Age': Avg Age of Houses in same city عمر متوسط خانه‌ها
- ۳) 'Avg. Area Number of Rooms': Avg Number of Rooms for Houses in same city متوسط تعداد اتاقهای هر خانه
- ۴) 'Avg. Area Number of Bedrooms': Avg Number of Bedrooms for Houses in same city متوسط تعداد اتاق خوابهای هر خانه
- ۵) 'Area Population': Population of city house is located in جمعیت شهر
- ۶) 'Price': Price that the house sold at قیمت فروخته شده هر خانه
- ۷) 'Address': Address for the house آدرس

۳. آیا ابهامی در تعاریف داده‌ها وجود دارد؟ خیر هیچ ابهامی وجود ندارد داده خالی یا missing value و مقدار تکراری هم ندارد و همه‌ی داده‌ها از یک نوع بودند و نیازی به تغییر نوع داده هم نبود
۴. آیا ممکن است در اندازه گیری متغیرها و یا ثبت داده‌ها خطایی وجود داشته باشد؟ بله امکان خطا هست ولی در داده‌های ما هیچ خطا یا missing value ای وجود ندارد چون همه ستونها دارای ۵۰۰۰ رکورد هستند.

۵. چه متغیرهای دیگری اگر وجود داشتند، میتواندست به حل مسئله کمک کند؟ اگر فقط مساحت را داشته باشیم از الگوریتم رگرسیون خطی استفاده میشود اگر چند متغیره باشد از الگوریتم polynomial regression (الگوریتم چند متغیره) استفاده میکنیم. در رگرسیون خطی مدل و شیب و خط برازش را تخمین میزنند. یکسری متغیر واقعی داریم و یکسری متغیری که پیش بینی شده Y^A و به ما در حل مسئله کمک میکند برای ارزیابی باید (متغیر پیش بینی شده - متغیر واقعی) را استفاده کنیم. متغیر X (مساحت) و متغیر Y (قیمت خانه است اینها متغیرهای واقعی هستند.

متغیرهای پیش بینی شده \hat{y} می باشد. مقدار β_1 شیب خط است و x مقدار مستقل است و y مقدار وابسته است. β_0 ثابت می باشد. و برای ارزیابی مقدار پیش بینی شده از فرمول $y - \hat{y}$ استفاده میشود (متغیر پیش بینی شده - متغیر واقعی). هرچه قدر فاصله کمتر باشد پیش بینی درست تر است

$$y = \beta_0 + \beta_1 x$$

۶. متغیرهای موجود از کدام نوعند (رسته ای - عددی)؟ عددی

۷. خلاصه آماری متغیرهای موجود چیست؟ برای تحلیل آماری از تابع `pairplot` و `distplot` استفاده میکنیم. برای اینکه ببینیم قیمت خانه نرمال هست یا خیر از تابع `distplot` استفاده کردیم. (خانه ها از یک میلیون تا دو میلیون هستند و بیشترین قیمت خانه روی ۱.۵ میلیارد است و این میانه دیتا است و این توزیع نرمالی ست) چولگی به چپ و راست ندارد - میانگین (`mean`) و انحراف معیار (`std`) دیتا فریم را بدست میاوریم.

۳) آماده سازی داده

۱. آیا نیاز به در آمیختن داده ها است؟ اقدامات و نتایج گزارش شود.

۲. آیا نیاز به پاکسازی داده است؟ اقدامات و نتایج گزارش شود. بله برای استاندارد کردن داده ها از `standardize` استفاده میکنیم تعداد اتاق خوابها ۳ تاست و جمعیت ۴۴ هزار تاست و داده ها نرمال نیستند و خیلی پایین و بالا هستند و باید داده ها استاندارد شوند. که از کتابخانه `sklearn.preprocessing` و از تابع `standardscaler` استفاده میکنیم و داده را استاندارد و یکسان میکنیم.

۳. آیا نیاز به تبدیل داده است؟ اقدامات و نتایج گزارش شود. فقط ستون `Adress` از نوع `object` است که انرا حذف میکنیم بقیه ستونها نوعشان `float ۶۴` است که نیازی به تبدیل ندارند. با دستور `info()` اطلاعات کاملی از نوع و خالی بودن سطرها نشان داده میشود.

۴. آیا نیاز به کاهش داده است؟ اقدامات و نتایج گزارش شود. خیر با استفاده از دستور `duplicate` داده های تکراری را پیدا میکنیم که در اینجا داده تکراری نداشت

۴) مدل سازی:

روی داده های آموزش، الگوریتم رگرسیون ساخته شود. اقدامات و نتایج گزارش شود. الگوریتم بانظارت را استفاده کردیم پس باید به ماشین مدل x, y را بدهیم تا یاد بگیرد الگویی را کشف کند از اینکه β_0, β_1 چند باشد

$$y = \beta_0 + \beta_1 x$$

۵) ارزیابی:

۶) مدل های ارائه شده، روی داده های آزمایش با استفاده از شاخص های ارزیابی رگرسیون خطی در یادگیری ماشین ارزیابی شوند. اقدامات و نتایج گزارش شود. انواع یادگیری ۱- نظارت شده (دیتاهایی که برچسب دارند و به مدل آموزش میدهند و مدل رابطه بین الگوها را کشف

میکند و براساس الگو خروجی جدیدی میدهد) ۲- بدون نظارت (هیچ لیبلای ندارد براساس دیتایی که به آن داده میشود خوشه هایی را مشخص میکند) ۳- تقویتی (براساس محیط است) که پروژه رگرسیون خطی براساس یادگیری نظارت شده پیاده سازی شده است. اگر خط برازش خوبی باشد پس مدل ما مدل خوبی است. میانگین خطای نهایی را بدست میآورند براساس آن نشان میدهند کدام روش مناسب است. برای ارزیابی روشهای متفاوتی داریم (متریک قدرمطلق خطا- متریک جذر خطا- جذر میانگین توان دو خطا) را داریم.

در مدل زدن باید متغیرهای وابسته و مستقل مشخص شوند با دستور iloc مشخص میکنیم کدام وابسته و کدام مستقل باشد. برای اینکه ببینیم مدل ما مدل خوبی هست یا نه آموزش و آزمایش train and test split تقسیم میکنیم. برای اینکه ببینیم مدل ما خوب کار کرده یک مقدار برای تست میزاریم. ۶۰ درصد آموزش و ۴۰ درصد تست. از کتابخانه sklearn استفاده میکنیم.

۷) ۲. چه پیشنهادهایی دارید تا نتایج در محیط واقعی، آزمایش گردد؟

۸) استقرار:

۹) حال اگر بخواهید چنین الگوریتمی را در مقیاس صنعتی توسعه دهید، به این فکر کنید با چه چالش هایی مواجه خواهید شد و

۱۰) برای آن چه راهکارهایی دارید. موارد زیر را گزارش کنید:

۱. چالشهای توسعه الگوریتم را بررسی کنید. ۱- نداشتن داده های کافی برای آموزش ۲- وجود داده های نویزی و بی

کیفیت ۳- بیش برازش ۴- کم برازش

۲. چه راهکارهایی برای حل آن ها دارید؟

۳. چه ملزوماتی برای ارائه آن راهکارها نیاز دارید؟

۱۱) نتیجه گیری:

۱. انجام این پروژه چه یادگیری برای شما داشت؟

۲. با چه چالشهایی مواجه شدید؟ چگونه آنها را حل کردید؟

۱۲) ارزیابی پروژه:

۱۳) کسانی که به بیش از ۶۰ درصد موارد بالا پاسخ صحیح داده باشند، نمره قبولی پروژه را دریافت خواهند کرد.

۱۴) پروژه را در سامانه تربیت مربی و گیت هاب خود بارگذاری کنید word ، html فایل

۱۵) لینک گیت هاب خود را هم در سامانه قرار دهید.