

به نام خدا

عنوان پروژه: دسته بندی مشتریان بر اساس رفتار آنها جهت ارائه خدمات

گردآورندگان: مهندس مونا رستگار، مهندس اشکان بویری

شرح پروژه:

- برای اینکه کار کردن با ستون ها راحت تر به وسیله پایتون انجام بشه نام آنها با معادل انگلیسی شان جایگزین شده
- نوع بعضی از ستون ها از فلو^۱ت^۲ ۶۴ به اینت^۳ ۶۴ تغییر داده شده.
- همبستگی بین ویژگی ها بدست آمده.
- مقادیر نامعلوم^۴ در میان داده ها وجود داشته که از داده ها حذف شدند.
- با رسم نمودار های متوجه شدیم که ویژگی ها زیر نه تنها تاثیر چندانی روی نتیجه نهایی مدل نخواهند داشت بلکه مدل رو دچار اشتباه خواهند کرد و هزینه محاسباتی اضافی را ایجاد می کنند با توجه به موارد ذکر شده این ستون ها از داده ها حذف شدند:
- میانگین فاصله بین روزهای انجام تراکنش و تعداد روز عدم فعالیت از آخرین تراکنش تا تاریخ روز و تعداد روز از افتتاح سپرده تا تاریخ آخرین تراکنش و ردیف و تعداد حساب و مانده در تاریخ ۹۸-۶-۳۱
- از ویژگی های خام دیتاست، ویژگی های جدید ، موثرتری استخراج شدند. با استفاده از تفاضل بین مقادیر ستون تعداد تراکنش بستانکار از ستون تعداد تراکنش بدهکار ستونی

^۱Float64

^۲Int64

^۳Nan

تحت نام دی بی ان سی دی^۴ ساختم که بر خلاف انتظارم کرولیشن بالایی با هیچ کدوم از ویژگی ها نداره اما این ستون را حذف نکردیم و ستون دیگه ای از تفاضل مقادیر میانگین تراکنش های بستانکار از ستون میانگین تراکنش بدهکار تحت نام دی اف ای سی دی^۵ ساختم این ستون کرولیشن مناسب با ستون میانگین موجودی (۹۵,۰) دارد، ویژگی های دیگری را میتوان از میان داده ها موجود استخراج کرد.

- از آنجاکه ما در این پروژه از الگوریتم های خوشه بندی استفاده خواهیم کرد و این الگوریتم ها به مقیاس داده ها حساسیت بالایی دارند باید داده هایمان در یک مقیاس باشند تا این الگوریتم ها عملکرد بهتری داشته باشند ما از هر دو روش استانداردسازی^۶ و نرمالیز کردن^۷ استفاده می کنیم یا استفاده از استانداردسازی مجموعه داده بانک_اس سی^۸ بدست آمده و با استفاده از نرمالیز مجموعه داده بانک_نرم بدست آمده که ما برای آموزش مدل ها از داده های بانک_نرم^۹ استفاده میکنیم.
- برای ارزیابی خوشه بندی های صورت گرفته با الگوریتم خوشه بند کی مینز و انتخاب تعداد خوشه مناسب از معیار اینرسی^{۱۰} و روش آرنج استفاده کردیم همچنین برای سنجش میزان اعتبار خوشه بندی با تعداد خوشه های مختلف در محدوده ۲ تا ۱۵ از معیار سایلیت اسکور^{۱۱} استفاده کرده ایم
- در گام بعدی برای اینکه دید بهتری از خوشه بندی های انجام شده بگیری از نمودار اسکتر پلات استفاده کردیم.
- سپس ارتباط بین ویژگی ها و خوشه ها را با معیار کرولیشن بدست آوردیم.

^۴DBNCD

^۵DFACD

^۶Standard Scaler

^۷Normalize

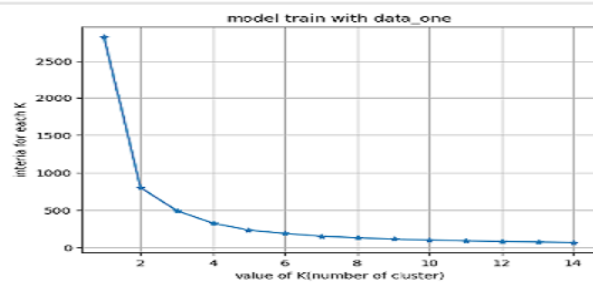
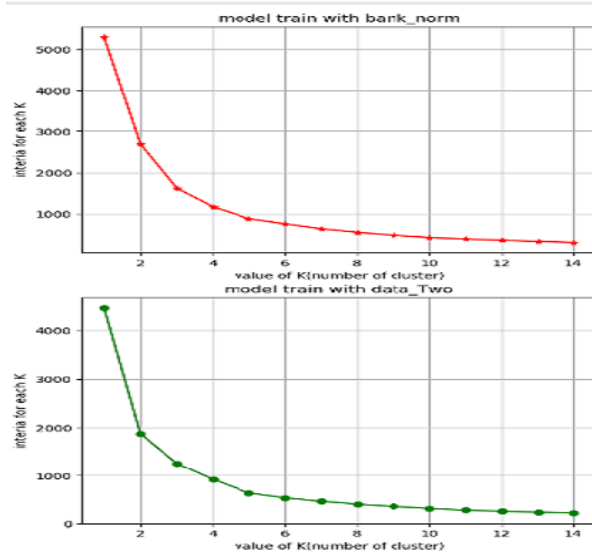
^۸Bank_sc

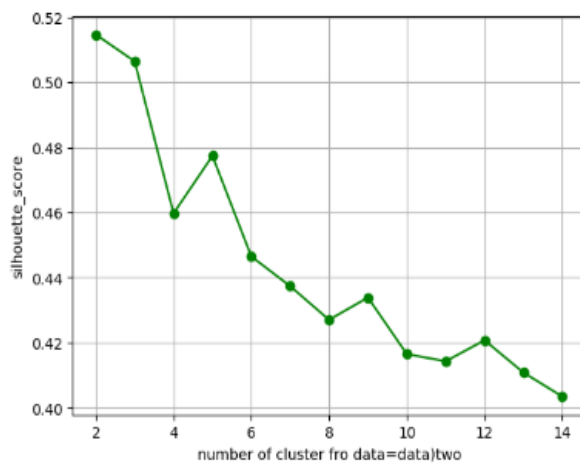
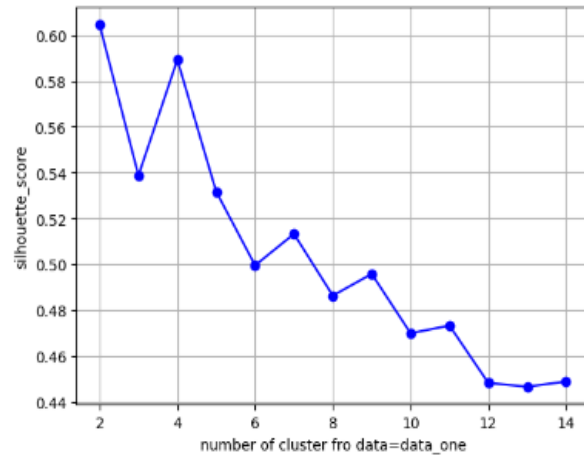
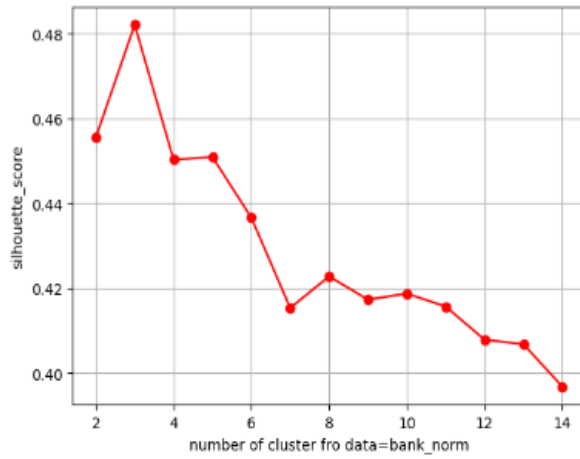
^۹Bank_norm

^{۱۰}Inertia

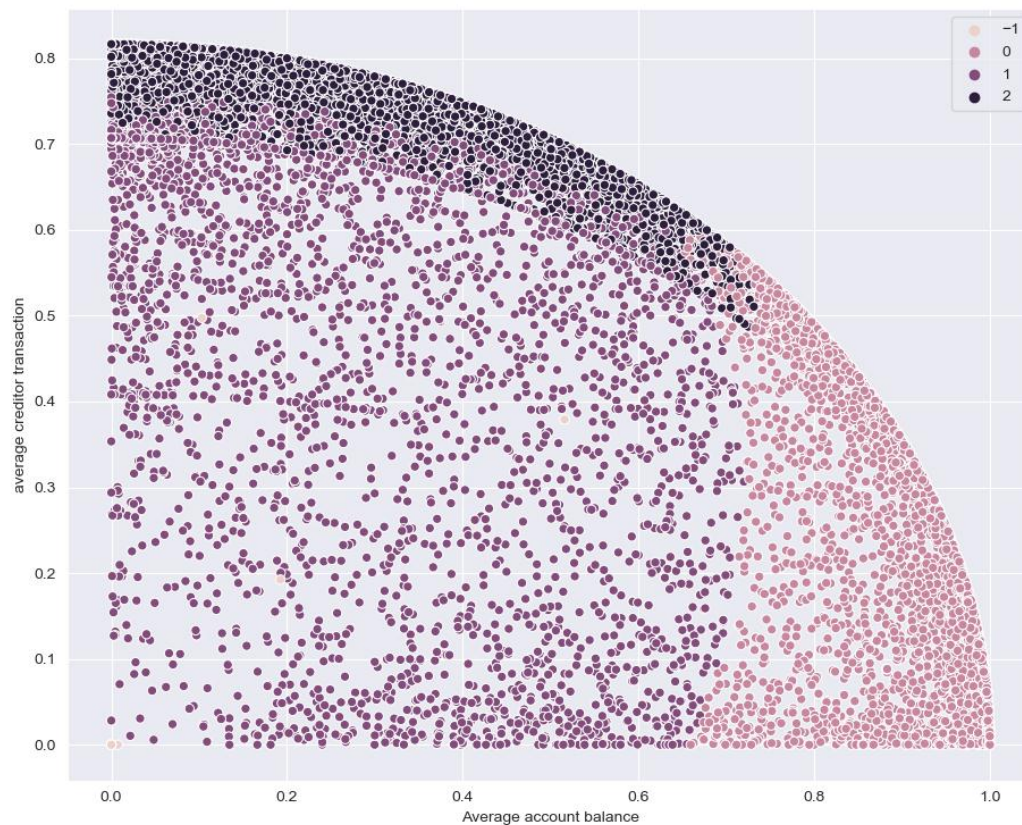
^{۱۱}Silhouette Score

- با تحلیل های صورت گرفته به این نتیجه رسیدیم که داده های موجود در خوشه دو نشان دهنده افراد خوش حساب بوده و داده های خوشه یک و صفر به ترتیب هر کدام نماین گر میان حساب و بد حساب می باشند.





- ۱۱ ردیف از داده ها با استفاده از الگوریتم دی بی اسکن نويز تشخيص داده شدند و از مجموعه داده ها حذف شدند.
- همچنین الگوریتم دی بی اسکن ترتیب خوشه ها را عوض کرده بدین صورت که داده ها در خوشه صفر، یک و دو به ترتیب بیان گر خوش حساب ها، میان حساب ها و بد حساب ها می باشند.



- در بخش دسته بندی از داده های بانک_نرم_دی بی استفاده کرده ایم، با استفاده از تکنیک های داده کاوی زیر مجموعه ای از داده ها ایجاد کرده و الگوریتم های دسته بندی مختلفی مانند جنگل تصادفی، ماشین بردار پشتیبان و.... را که پارامترهای این الگوریتم ها با استفاده از تابع گرید سرچ سی وی بهینه شده است با استفاده از داده های انتخاب شده آموزش دادیم. مدل های توسعه داده شده همگی دقت و عملکرد بالایی داشته.