



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده ریاضی و علوم کامپیوتر
کارشناسی ارشد علوم کامپیوتر گرایش داده کاوی
پروژه شماره ۴ درس یادگیری ماشین

نگارش

حدیث حق شناس جزی

استاد راهنما

محمد اکبری

استاد مشاور

محمد علی سفیدی اصفهانی

دی ۱۴۰۱

چکیده

در این پروژه پس از فراخوانی داده ۱۶۷ کشور (شامل ۱۰ ستون از اطلاعات سلامت و درآمد و میزان امید به زندگی و ...) و انجام بررسی های ابتدایی به کمک معیار های elbow method و silhouette score از چند طریق با بررسی میزان نزدیکی این عدد به ۱ و یا تعداد اجزای منفی و ... به انتخاب تعداد خوشه بندی پرداخته ایم سپس با انواع الگوریتم های kmeans, spectral , GMM خوشه بندی را بر روی داده انجام داده و به کمک اضافه کردن لیبل هر خوشه که به داده اضافه شده است، همه ستون ها را نسبت به یکدیگر نمودار نقطه ای کشیده ایم و هر رنگ را به یک خوشه اختصاص داده ایم.

قابل ذکر است که silhouette score در ادامه توضیح داده میشود همچنین نمودار های آن نیز در کد ها قابل مشاهده است. به طور کلی هرچه مقدار این عدد به ۱ نزدیکتر باشد یعنی خوشه بندی ما دارای تراکم بیشتر و همچنین جدا پذیری بهتر است و هرچه به منفی ۱ نزدیک شود برعکس این امر اتفاق می افتد. اما باید توجه داشت که میزان این عدد از ۰ نباید کمتر شود (در این صورت خوشه بندی به درستی انجام نمیپذیرد) لذا ما به دنبال عددی برای خوشه بندی هستیم که تعداد مولفه منفی آن صفر یا کمتر از تعداد خوشه بندی دیگر باشد و همچنین عدد بالاتری نسبت به اعداد دیگر نسبت داده شده به خوشه بندی های دیگر داشته باشد.

توضیحات کلی :

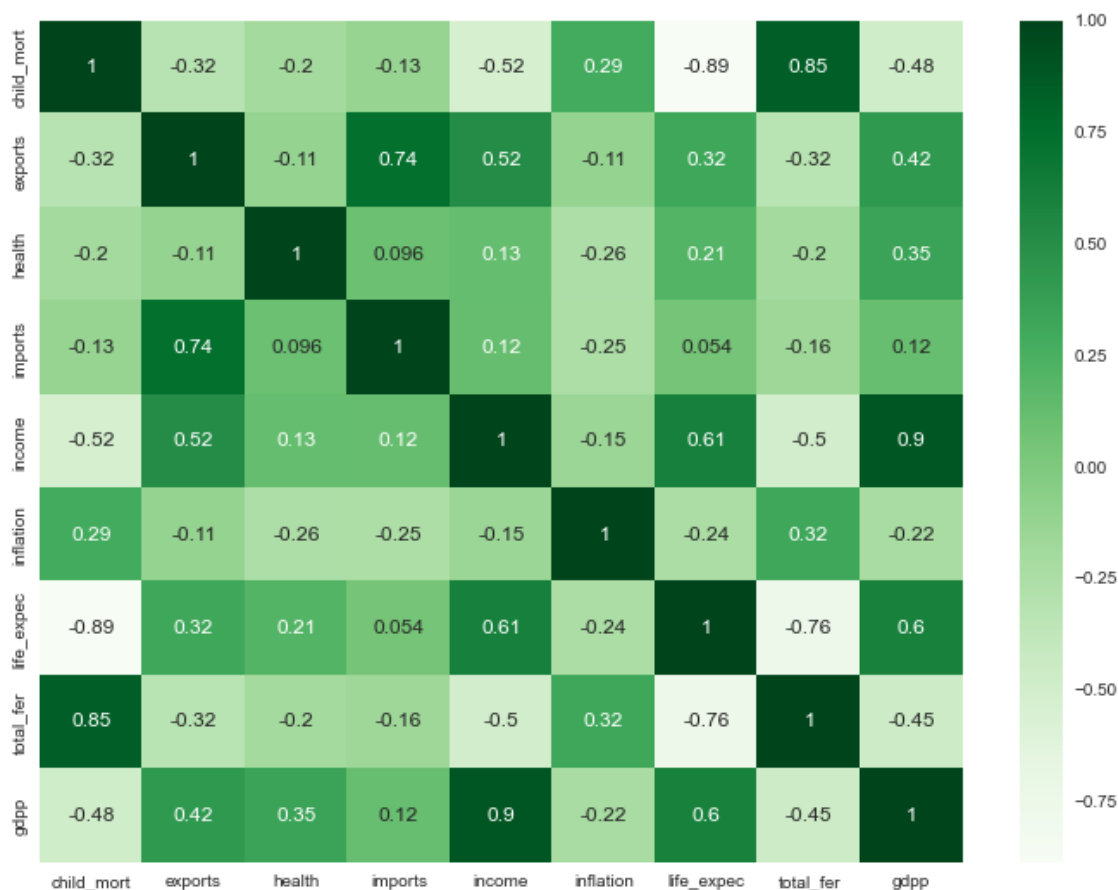
- تمرین تشریحی در یک فایل pdf موجود میباشد.
- توضیح تمامی قسمت کد ها در فایل ژوپیتر قابل مشاهده است. (متاسفانه ذخیره به صورت تکست کار نکرد و امکان آوردن کد ها در این فایل نبود)
- به جای دو به دو نمودار کشیدن در بخش انتخاب ۳ مولفه، تمامی مولفه ها نسبت به یکدیگر نمودار شده اند (به نظر بهتر می آمد)
- پس از مرحله اعمال PCA و ایجاد یک دیتا فریم جدید، این دیتا فریم در فایلی به نام ml4.pca آمده و مراحل خواسته شده بر روی این داده در یک فایل جدا انجام گرفته است .

معیار silhouette score :

یکی از روش های ارزیابی خوشه بندی، معیار (Silhouette) است (که به آن معیار نیم رخ نیز گفته میشود). این معیار هم به پیوستگی درون خوشه ها و هم به میزان تفکیک پذیری آن ها بستگی دارد. مقدار نیم رخ برای هر نقطه، میزان تعلق آن را به خوشه اش در مقایسه با خوشه مجاور اندازه می گیرد. مقدار این شاخص بین ۱- تا ۱+ تغییر می کند. مقدار نزدیک به ۱ بیانگر انطباق خوب بین نقطه و خوشه اش نسبت به خوشه مجاور است. اگر معیار نیم رخ برای همه نقاط درون خوشه ها نزدیک به ۱ باشد، عمل خوشه بندی به درستی انجام شده است. در حالیکه کوچک بودن مقدار نیم رخ برای خوشه ها، بیانگر ضعیف بودن نتایج خوشه بندی است که ممکن است به علت انتخاب نامناسب تعداد خوشه ها (k) نیز باشد. اگر میانگین مقدار نیم رخ برای نقطه های هر خوشه را محاسبه کنیم، معیاری برای ارزیابی هر خوشه بدست می آید. همچنین میانگین کل مقدارهای نیم رخ نیز معیاری برای ارزیابی عملیات خوشه بندی محسوب می شود.

برای تفسیر این معیار، از نموداری استفاده می شود که میزان انطباق هر نقطه را با خوشه خودش نمایش می دهد. در تصویر زیر این نمودار دیده می شود. محور افقی نقطه ها و ستون ها، مقدار معیار نیم رخ برای آن نقطه است. همچنین میانگین شاخص نیم رخ برای همه نقاط نیز در نمودار مشخص می شود. اگر در این نمودار بعضی نقاط دارای مقدار نیم رخ منفی باشند نشان می دهد که ممکن است به درستی خوشه بندی نشده باشند و به خوشه مجاور تعلق داشته باشند.

توضیح نمودار ها:



در نمودار بالا مشاهده میشود که بین ستون های gdp و income ضریب همبستگی ۰,۹ وجود دارد . همچنین بین دو ستون child_mort و total_fer این مقدار ۰,۸۵ است. هر چه ضریب همبستگی بیشتر باشد، امکان پیش بینی مقدار یکی از متغیرها برحسب دیگری بیشتر است لذا دو ستون از بین این ۴ ستون حذف میکنیم چون اطلاعات نزدیک به یکدیگر در بر دارند.

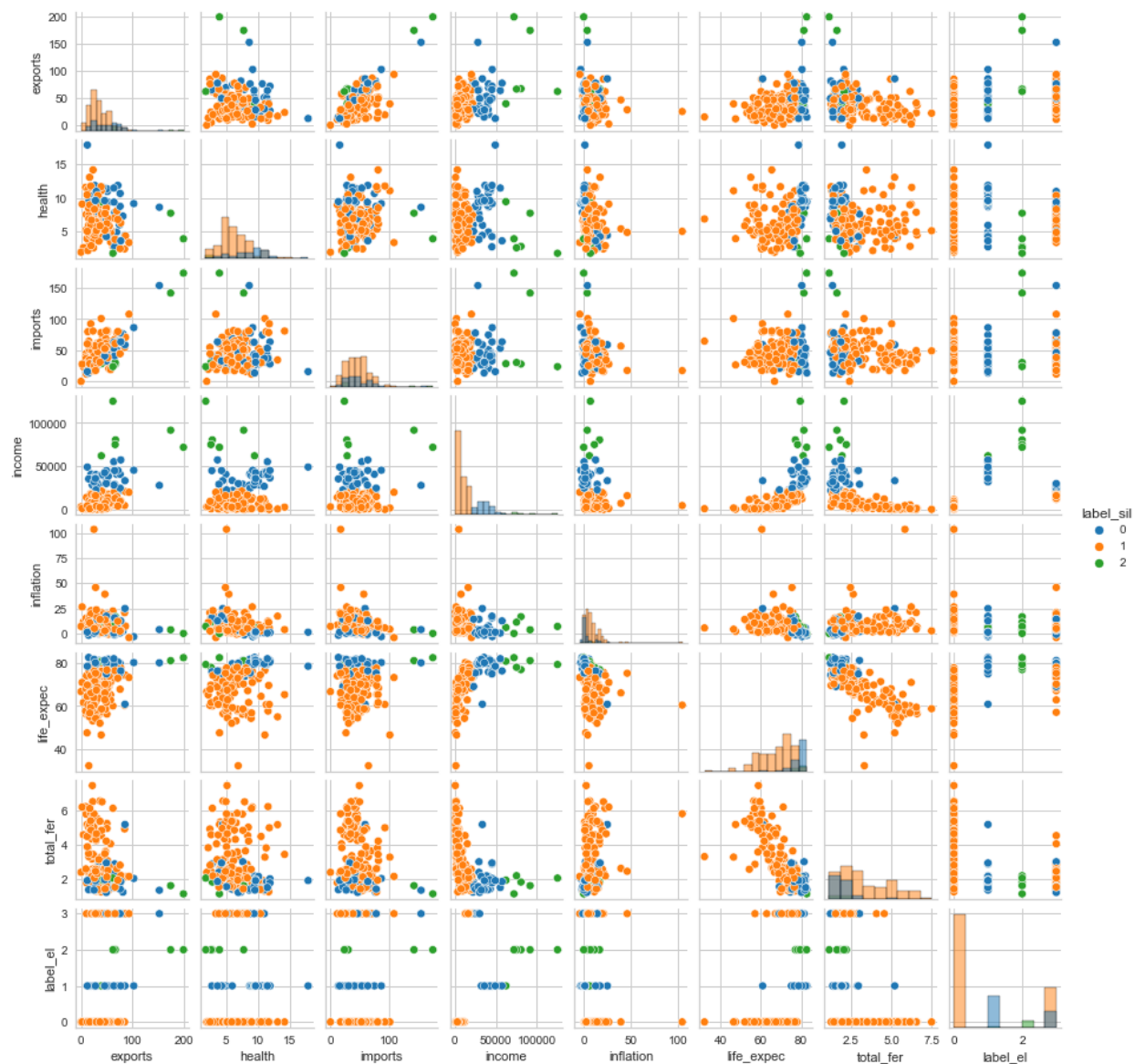


خوشه بندی kmeans با تعداد خوشه elbow method:

نسبت `income` به `life_expect`: داده ها به صورت `curve` هستند و خوشه بندی به خوبی داده ها را از یکدیگر تمیز داده است. میتوان گفت که داده های پرت در خوشه سبز قرار گرفته اند. هرچه امید به زندگی بیشتر میشود، میزان درآمد نیز افزایش دارد مخصوصاً در نقاط انتهایی که با شیب زیاد میزان امید به زندگی در کشور هایی با درآمد بالا افزایش یافته است (و بالعکس)

نسبت `income` به `imports`: همبستگی داده ها خنثی است و داده ها به خوبی خوشه بندی شده اند. خوشه سبز رنگ شامل داده های پرت میشود. با احتمال نسبتاً ضعیف به میزان درآمد، واردات وجود داشته است.

نسبت income به health : لزوماً با افزایش درآمد کشور وضعیت سلامت بهبود نیافته است اما گروه هایی که در یک وضعیت نسبی سلامت هستند در خوشه های یکسان قرار گرفته اند.

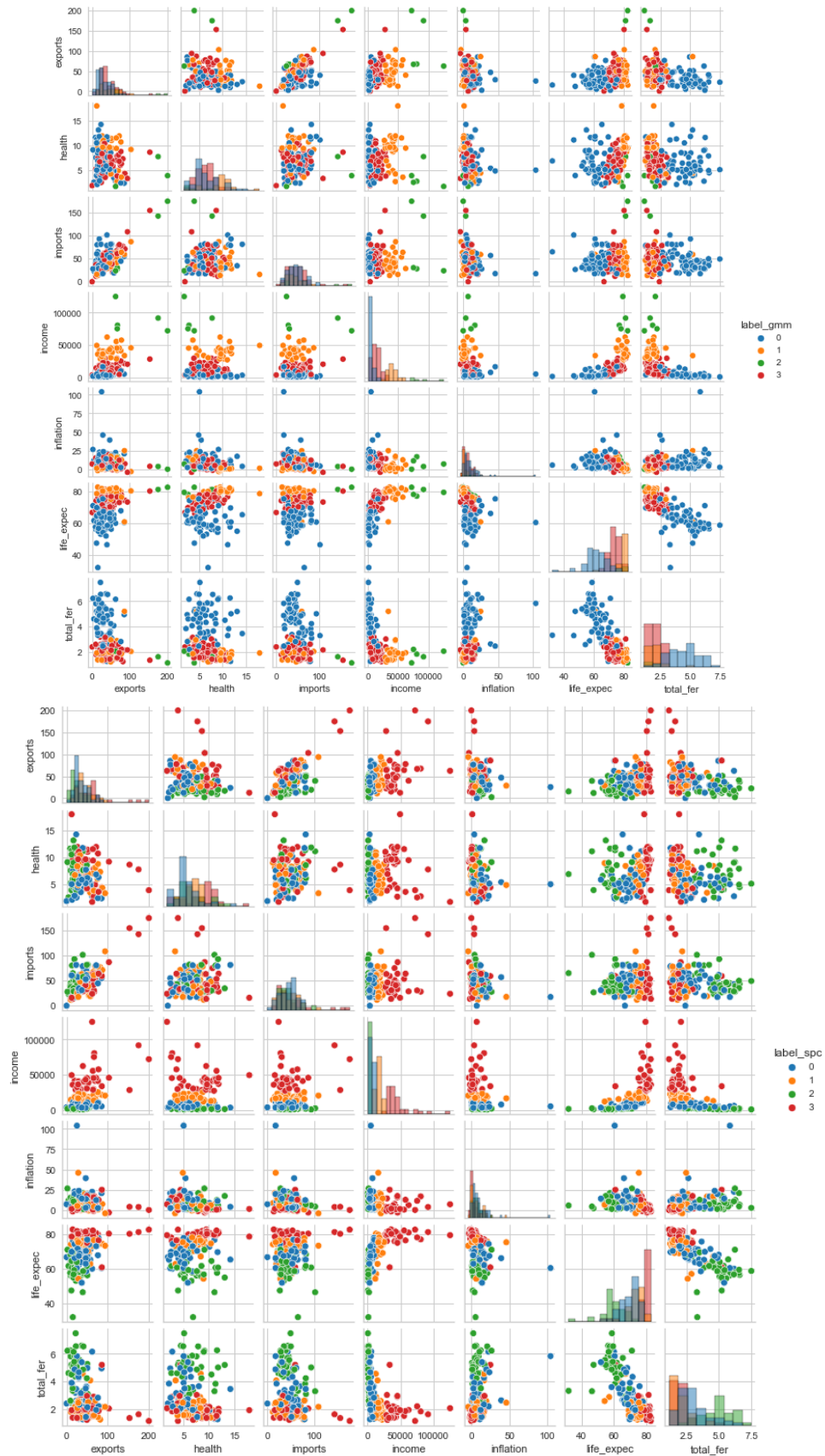


خوشه بندی kmeans با silhouette score :

نسبت income به life_expect : داده ها به صورت curve هستند و خوشه بندی به خوبی داده ها را از یکدیگر تمیز داده است . میتوان گفت که داده های پرت در خوشه سبز قرار گرفته اند. هرچه امید به زندگی بیشتر میشود، میزان درآمد نیز افزایش دارد مخصوصاً در نقاط انتهایی که با شیب زیاد میزان امید به زندگی در کشور هایی با درآمد بالا افزایش یافته است (و بالعکس)

نسبت income به imports : همبستگی داده ها خنثی است و داده ها به خوبی خوشه بندی شده اند . خوشه سبز رنگ شامل داده های پرت میشود. با احتمال نسبتا ضعیف به میزان درآمد, واردات وجود داشته است.

نسبت income به health: مانند خوشه بندی قبلی منتها به طور کلی به نظر من خوشه بندی ۴ تایی بهتر از خوشه بندی ۲ تایی عمل کرده است و elbow method عدد بهتری به ما داده است .

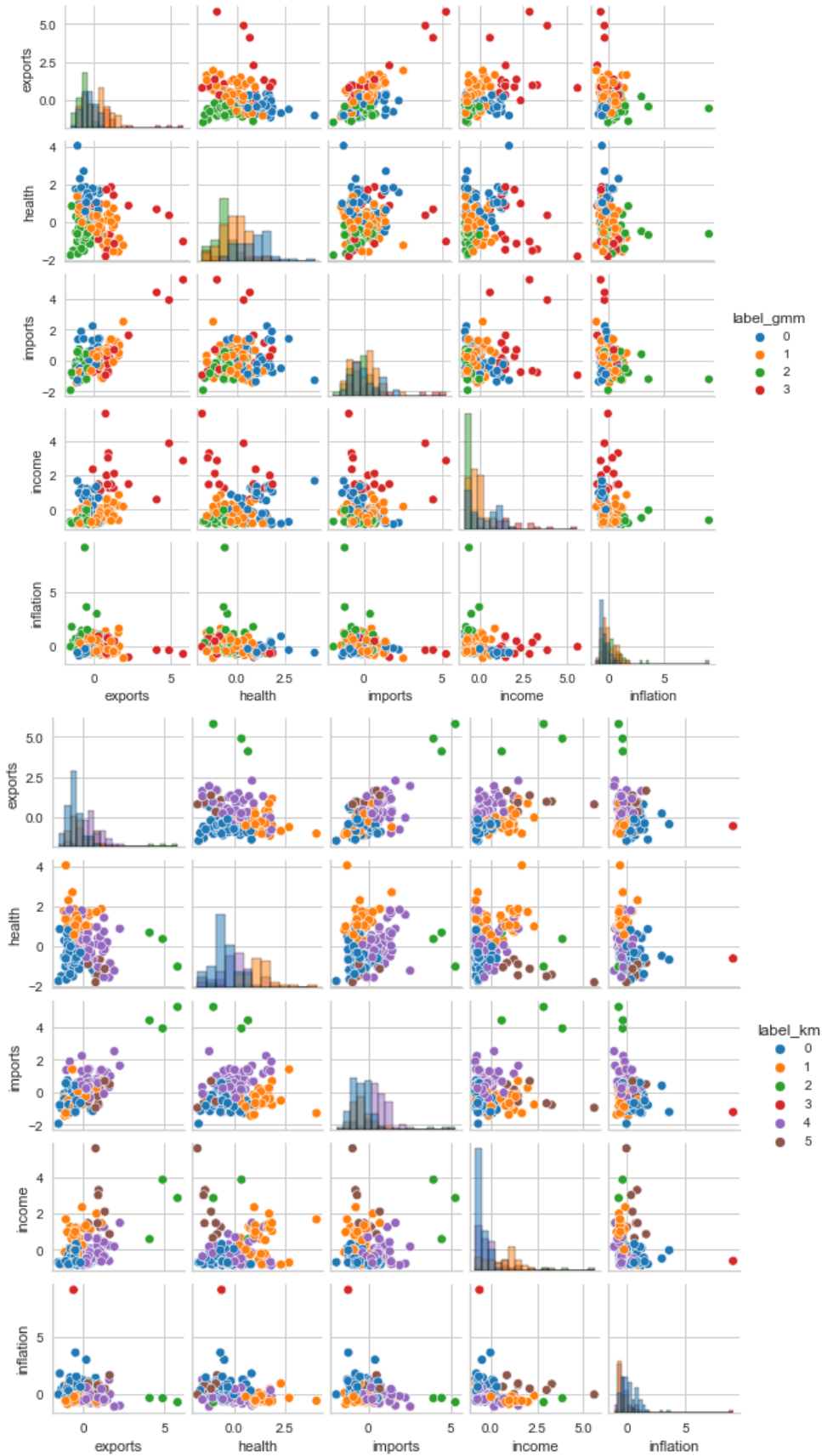


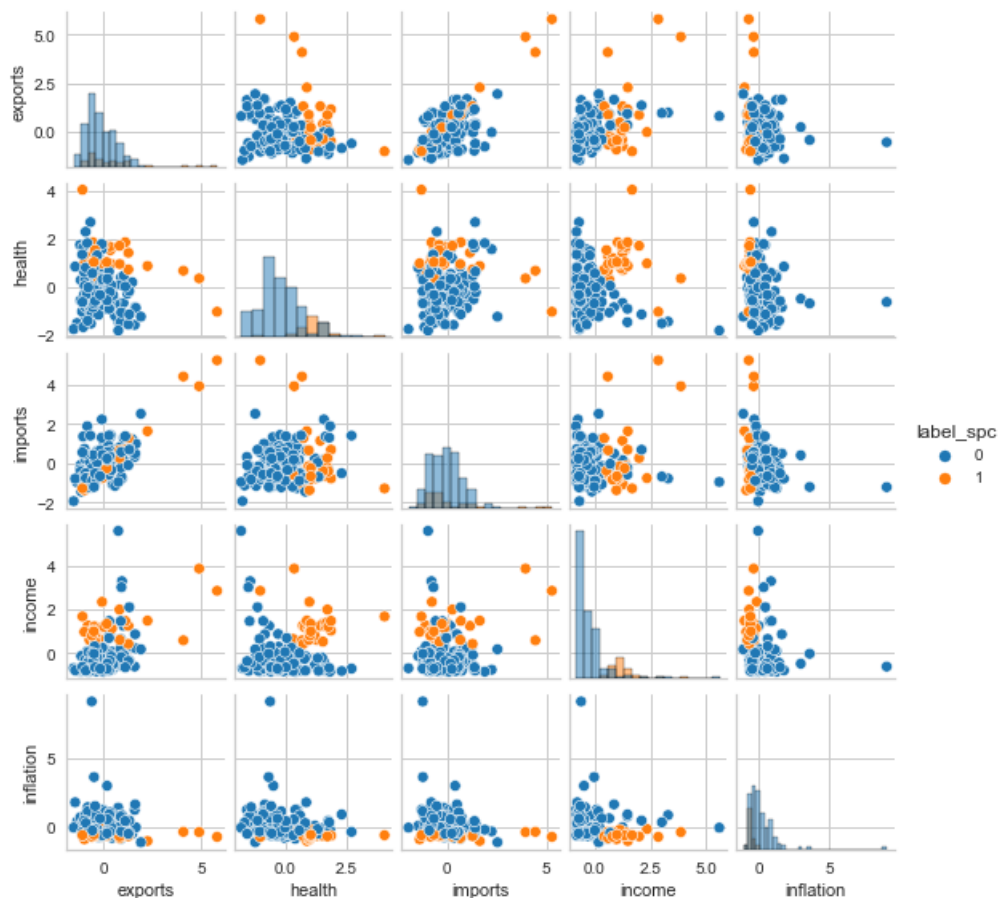
نمودار های GMM و Spectral clustering :

نسبت income به life_expec : نتایج مانند نمودار پیشین است . به نظر می آید که در GMM خوشه بندی کمی دقیق تر صورت گرفته است

نسبت income به imports : نتایج مانند نمودار پیشین است . به نظر می آید که در spectral خوشه بندی کمی دقیق تر صورت گرفته است

نسبت income به health : نتایج مانند نمودار پیشین است . به نظر می آید که در spectral خوشه بندی کمی دقیق تر صورت گرفته است



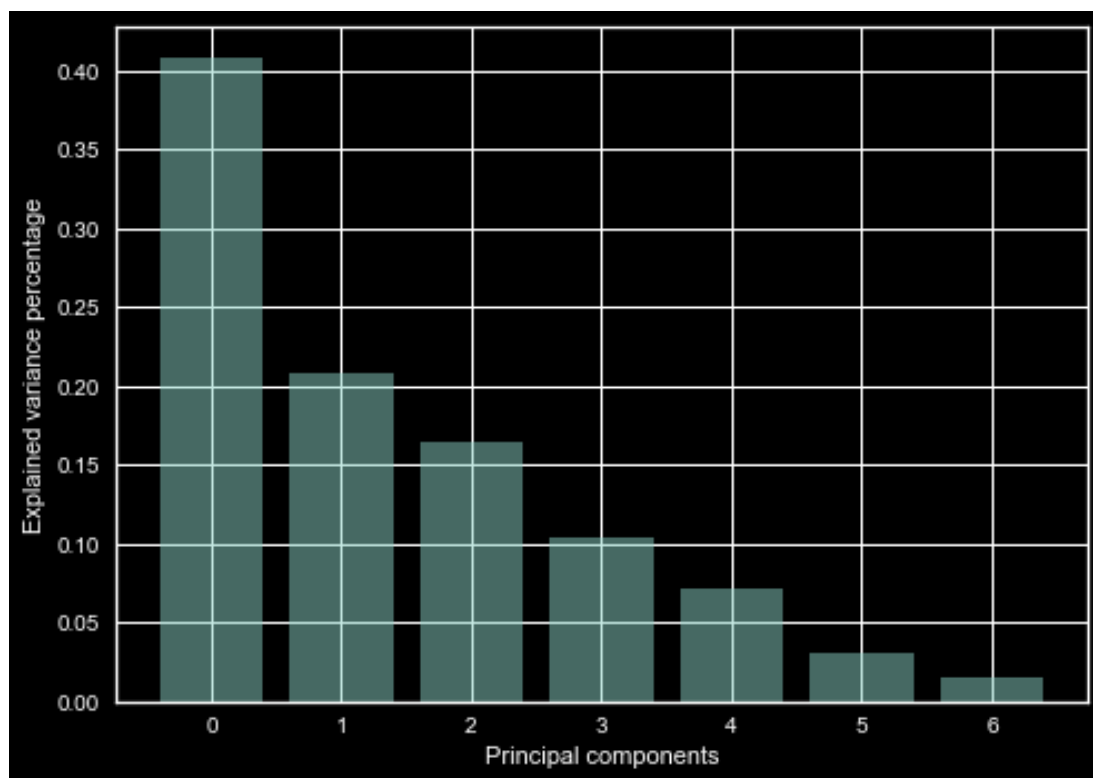


به ترتیب بعد از اعمال pca سه نمودار از دسته های kmean و GMM و spectral clustering :

نسبت income به infilation : داده ها به صورت خطی هستند و خوشه بندی کیفیت نسبتا پایین تری دارد . هرچه میزان درآمد کشور بالا رفته است تاثیری بر inflation نداشته است. همچنین در خوشه بندی spectral بهتر از دوتای دیگر خوشه بندی صورت پذیرفته است.

نسبت health به imports : نسبتا هرچه میزان سلامت بالاتر رفته است میزان import هم افزایش داشته است و بهترین خوشه بندی به کمک الگوریتم GMM انجام شده است.

نسبت health به income : بهترین خوشه بندی توسط الگوریتم های spectral و kmeans صورت گرفته است. در خوشه بندی آخر میتوان به خوبی مشاهده کرد که داده های پرت در یک خوشه قرار گرفته اند .



(نمودار مربوط به pca . همانطور که مشاهده میشود ستون های اول تا پنجم بیش از ۹۰٪ تاثیر گذاری و پراکندگی را در دیتا دارند پس حفظ میشوند و مابقی را حذف میکنیم.

پایان