

توجه: برخی توضیحات ابتدایی در نوت بوک داده شده است.

## درخت تصمیم [decision-tree-and-more/jupyter-titanic.ipynb]

### مرحله اول | بررسی داده‌ها و انتخاب فیچرها و تبدیل فیچرهای غیرعددی به عددی

#### آنالیز داده‌ها و رسم نمودار

با توجه به نمودار رسم شده در [5] In با چشم مشخص می‌شود که زن بودن، در بقا خیلی موثر است؛ پس قطعا فیچر sex برای درخت تصمیم انتخاب می‌شود.

در [6] In، نمودار مربوط به مقایسه با توجه به سنین مختلف رسم شده که در همه‌ی رنج‌های سنی نزدیک پنجاه درصد است. به جز در حدود 20 تا 30 سالگی و پیرمردان و زنان لب گور. همانطور دیده می‌شود (در نوت بوک) با شرط پیر بودن، درخت تصمیم می‌تواند تصمیم قاطعی بگیرند.

این دو ویژگی بر اساس "اول زن‌ها و بچه‌ها" بررسی شدند. ویژگی بعدی که در [7] In نمودار آن رسم شده، بررسی می‌کند که آیا کسی که پول بیشتری داده، لیاقت بیشتری برای زندگی دارد یا نه؟ در این نمودار تجمعی شیب خط مربوط به کشته شدگان در ابتدا خیلی زیاد است، این نشان می‌دهد که افراد زیادی که کشته شده‌اند (بیش از چهارصد نفر)، پول خیلی کمی برای بلیط داده‌اند. علاوه بر این کسانی که زنده مانده‌اند مجموع پولشان به پانصد رسیده و کسانی که کشته شده‌اند، به سیصد هم نرسیده، این درحالی است که تعداد زنده ماندگان حدود پانصد و کشته شدگان، سیصد است.

#### انتخاب ویژگی‌ها

در مرحله‌ی بعد در [8] In با ماتریس وابستگی بررسی کردم که آیا Pclass و Fare وابستگی زیادی دارند، تا یکی از آن‌ها را حذف کنم. نتیجه -0.57 بود، که به نظر من کافی نیست، پس هیچ فیچری در این مرحله حذف نشد.

درباره‌ی سایر ویژگی‌ها در نوت بوک توضیحات داده شده است.

#### پیش‌پردازش

در مرحله‌ی بعدی، شروع به حذف داده‌های NaN کردم. برای سن از میانگین استفاده کردم و برای Cabin از No room.

همچنین برای Embarked با کمک mode گرفتن از این ستون، مقادیر خالی، جایگزین شده‌اند. در نهایت در test.csv یک سطر از Fare خالی بود، که با میانگین مرتبط، جایگزین شد.

در [13] In شروع به تبدیل داده‌های عددی به غیر عددی کردم که به شکل زیر می‌باشد:

متغیرهای جدید	متغیرهای قبلی	فیچر
-1	Female	Sex
1	male	
0	S	Embarked
1	C	
2	Q	

دو ویژگی SibSp و Parch مربوط به خانواده هستند و می‌توانند با هم جمع شوند، برای این یک ویژگی FamSize تعریف شده که مجموع این دو فیچر و یک (خود طرف) است.

در نهایت با دسته‌بندی Cabin با توجه به حرف اولشان، می‌توان اطلاعات نسبی بدست آورد این مرحله به پایان می‌رسد. این دسته‌بندی در [16] انجام شده که در آن کابین‌ها با شروع A تا G جایگزین 1 تا 7 می‌شوند و No room ها با صفر جایگزین می‌شوند. اسم این فیچر جدید cCabin است.

نتایج نهایی در [16] Out قابل مشاهده است.

## مرحله دوم | آموزش، تست و توییک

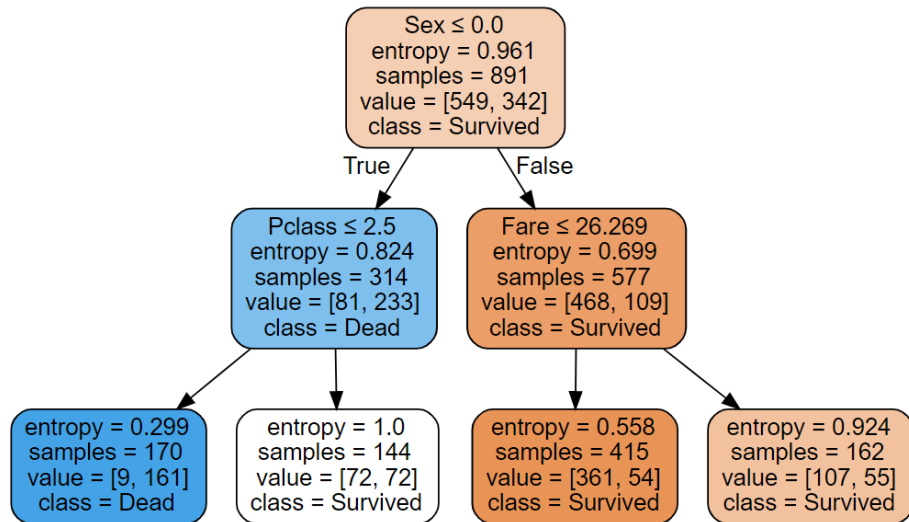
### آموزش

در اینجا، داده‌ها را با نسبت ۸۰ به ۲۰ به داده‌های آموزش و تست اختصاص دادیم، تا بتوانیم با توجه به نتایج، خروجی نهایی را بهبود بدهیم. این کار با کمک shuffle انجام شده است.

فیچرهای انتخاب شده بصورت زیر می‌باشند

- Age
- Sex
- Fare
- Pclass
- Embarked
- FamSize
- cCabin

برنامه با اجرا، درخت را نمایش می‌دهد و همینطور فایل پی دی اف آن را خروجی می‌دهد. برای نمونه درخت با عمق ۲ و با استفاده از انتروپی به شکل زیر است



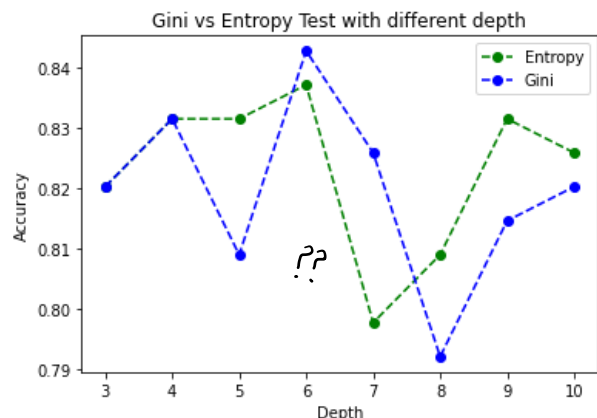
با تغییرات عمق و معیار تقسیم نتایج زیر حاصل شد

عمق	روش	میانگین آموزش	میانگین تست
۴	Entropy	۸۲	۸۲
۵	Gini	۸۳	۸۲
	Entropy	۸۳	۸۰
۶	Entropy	۸۴	۸۳
	Gini	۸۵	۸۲

در نهایت با بررسی میانگین‌ها، عمق ۶ با روش gini انتخاب شد. در حالت‌های خاصی، تست نتیجه‌ی زیر ۸۰ هم می‌دهد ولی در حالت کلی، بیش از ۸۰ است و همیشه بیش از ۷۵ است. تست‌های Kaggle هم نتوانست از نمره ۰.۷۸ بیشتر شود.

در نهایت در [24] In تاثیر عمق در معیارهای مختلف بررسی شده (که نسبتاً رندوم هستند) و پس از آن تخمین نهایی از test.csv در فایل output.csv ذخیره می‌شود.

این نمودار با shuffle یکسان، random state یکسان و تغییرهای یکسان رسم شده.



## جنگل تصادفی [decision-tree-and-more/jupyter-titanic.ipynb]

در [15] In برای بهبود نتایج، کلاس شخصیتی را به فیچرها اضافه کردم. این کار تاثیر نسبی در نتایج نهایی داشت. مقایسه‌هایی که در این تمرین انجام شده با تمرین قبلی است و برای این مقایسه‌ها [15] In را اجرا نکردم تا مقایسه‌ها در شرایط یکسان باشد.

در [24] In جنگل تصادفی بهینه را با روش gini و با 500 کلاس‌بند تشکیل دادم. مقدار  $n\_jobs=-1$  تعیین می‌کند که عملیات بصورت موازی انجام شود تا به سرعت بیشتری برسیم. ساده بودن درخت‌های تصمیم موجود در این جنگل با  $min\_samples\_split = 10$  حاصل شده که باعث می‌شود نودها با کمتر از 10 داده را بیشتر از این تقسیم نکنیم.

[25] In و [26] In به ترتیب دقت داده‌های آموزش و تست را چاپ می‌کنند.

در تمرین قبل، بهترین پاسخ برای درخت تصمیم، در حالت آموزش دقت ۸۷ و در حالت تست دقت ۸۱ بدست آمد. این مقادیر با پارامتر عمق ۶ و روش gini بدست آمدند. همانطور که در تمرین قبل ذکر شد، در حالت‌های خاصی، درخت تصمیم نتیجه‌ی زیر ۸۰ می‌دهد ولی هیچوقت کمتر از ۷۵ نیست و میانگین ۸۰ دارد.

در این تمرین بهترین نتیجه‌ی جنگل تصادفی، در حالت آموزش ۸۹ و در حالت تست ۸۴ شد. این حالت با ۵۰۰ کلاس‌بند با روش gini و عمق متغیر بر اساس تعداد اعضای هر نود، بدست آمد. میانگین جنگل تصادفی برای آموزش ۸۷ و برای تست ۸۳ بدست آمد.

نمودار رسم شده در [27] Out نشان‌دهنده‌ی اهمیت فیچرها در جداسازی داده‌های هستند.

جدول رسم شده در [28] Out مقایسه‌ی زمان انجام را برای درخت تصمیم و جنگل تصادفی انجام می‌دهد. برای مقایسه‌ی سرعت، با استفاده از کتابخانه زمان پایتون، به نتیجه‌ی زیر رسیدم:

	زمان آموزش میلی‌ثانیه	زمان تست میلی‌ثانیه
درخت تصمیم	4.00	0.09
جنگل تصادفی	444.20	104.13

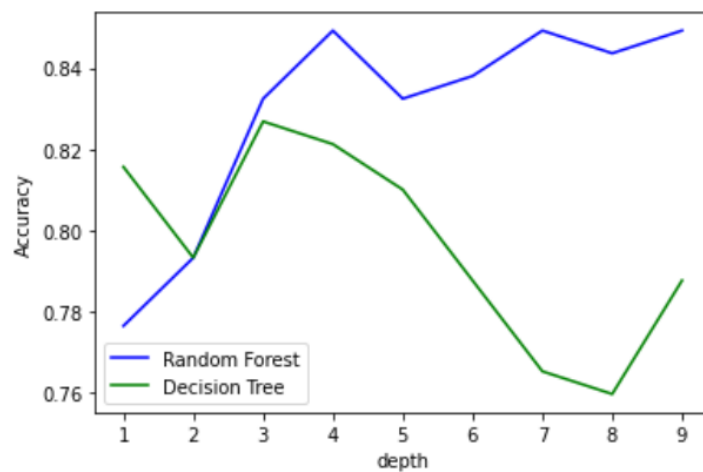
زمان‌های بدست آمده، زمانی است که در آن بهترین نتیجه‌ی درخت تصمیم و جنگل تصادفی بدست آمد.

با تغییر مقادیر عمق و روش جداسازی در [30] In، جدول زیر میانگین دقت‌ها را نشان می‌دهد.

عمق	روش	آموزش	تست
3	Gini	83	80
	Entropy	83	79
6	Gini	86	84

	entropy	86	84
--	---------	----	----

برای مقایسه‌ی کلی بین این دو روش، دقت تست را بر اساس عمق‌های مختلف تست کردم و به نتیجه‌ی زیر رسیدم:



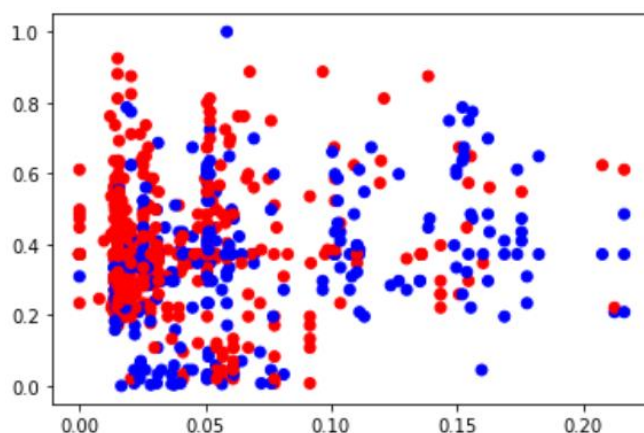
همانطور که مشخص است، به طور کلی جنگل تصادفی با صد کلاس‌بند از یک درخت تصمیم، بهتر عمل می‌کند.

## ماشین بردار پشتیبان [decision-tree-and-more/jupyter-titanic.ipynb]

برای پیش‌پردازش در [31] In با استفاده از نرمال‌سازی کتابخانه sklearn، داده‌ها را برای آموزش آماده کردم. در اینجا از روش minmax استفاده شده است.

مدل با کرنل خطی در [32] In آموزش داده شده است. دقت میانگین برای داده‌های آموزشی برابر 78 درصد بود و برای داده‌های تست به دقت 79 رسیدم. همانطور که مشخص است، SVM از روش جنگل تصادفی و درخت تصمیم بدتر عمل کرده است.

در [34] In اسکترپلات براساس age و fare رسم شده و مشخص است که با توجه به اینکه ابعاد بالا می‌باشد و پیچیدگی خطی نیست، نمی‌توان داده‌ها را بصورت خطی از هم جدا کرد و به همین دلیل است که کرنل خطی جواب نمی‌دهد و نیاز به جداساز پیچیده‌تری داریم.



در [33] In با کمک کرنل rbf و تنظیماتی مثل  $C=1$  و  $\gamma='scale'$  به نتایج به مراتب بهتری رسیدیم. به طور میانگین SVM با این کرنل به میانگین دقت آموزشی 81 و تست 83 می‌رسیم که از کرنل خطی بهتر است و تقریباً به دقت جنگل تصادفی می‌رسد.

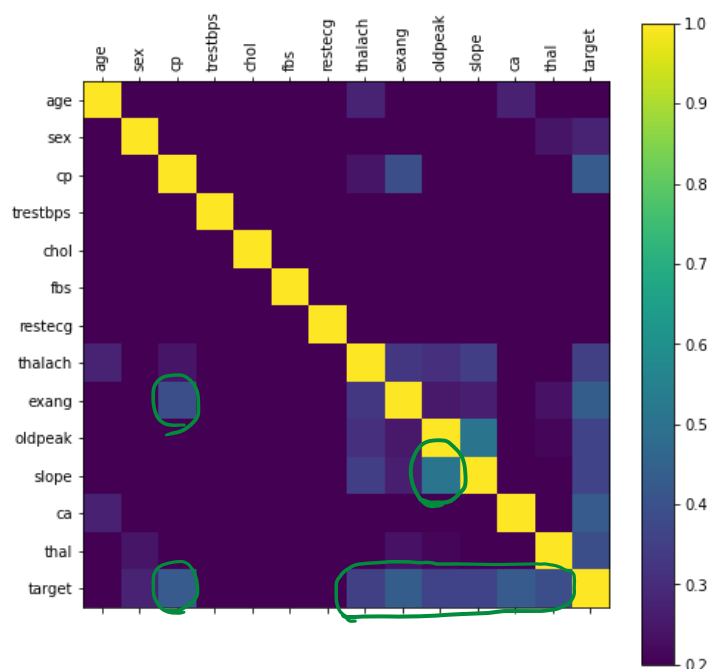
دلیل برتری کرنل rbf این است که پیچیده‌تر است و در واقع داده‌ها بصورت خطی قابل جدا شدن نیستند. به عبارت بهتر مدل با کرنل خطی bias دارد و underfit می‌شود ولی rbf این مشکل را ندارد و با انعطاف بیشتری می‌تواند جداسازهای بهتری بسازد.

## مرحله اول | تحلیل و آماده سازی

### تحلیل

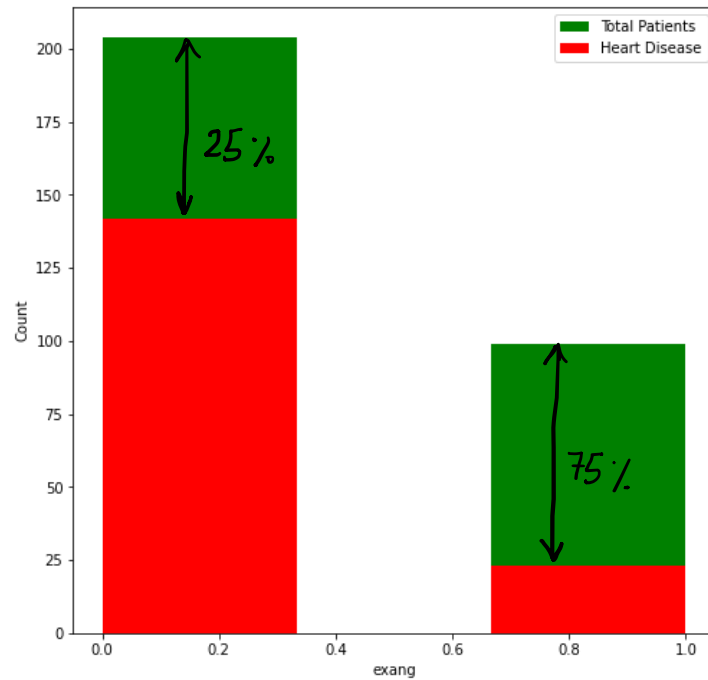
در ابتدا بررسی شده که تعداد داده‌ها در کلاس‌های مختلف، متناسب باشند تا در Naïve Bayes به مشکل برنخوریم. اینکار در [7] In انجام شده و نتیجه نشان می‌دهد که توزیع داده‌ها نسبتاً یکسان است.

در [9] In ماتریس وابستگی فیچرها رسم شده است. در اینجا برای وضوح، وابستگی‌های کوچکتر است 0.2 را هم‌رنگ صفر در نظر گرفته ام.

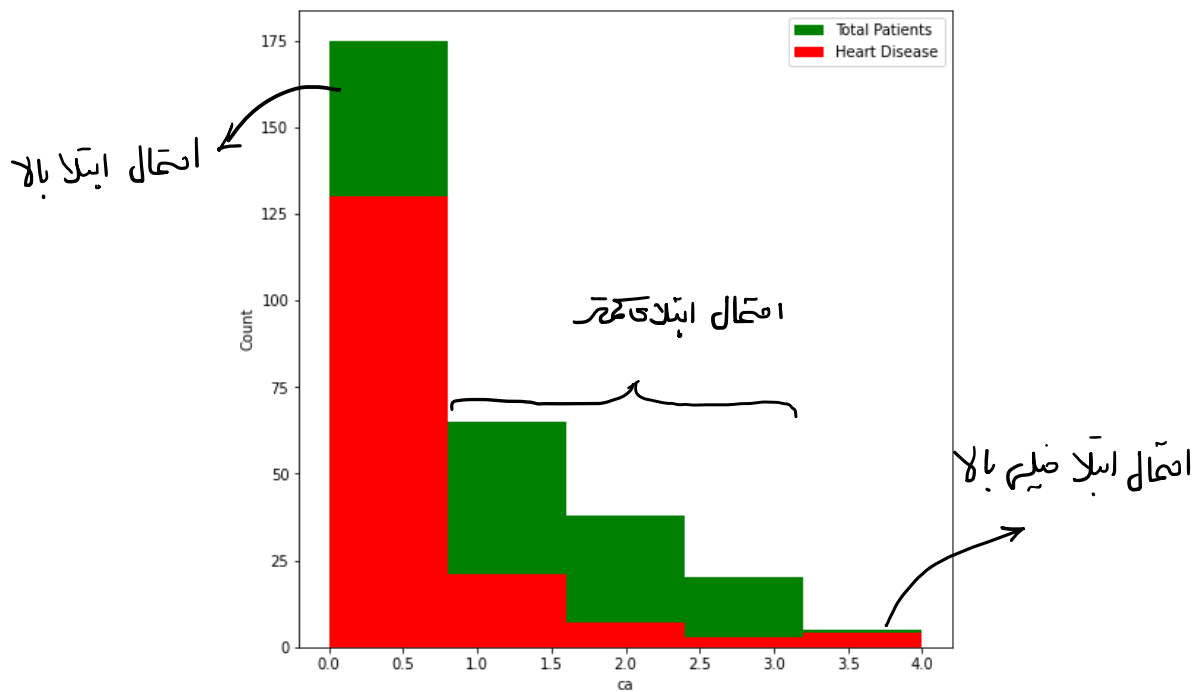


همانطور که در گوشه مشخص است، thal، ca، exang و cp با target ارتباط بیشتری دارند. در وسط تصویر هم slope با oldpeak وابستگی دارد که می‌توانیم یکی از آنها را حذف کنیم. در اینجا بدلیل محدود بودن داده‌های آموزشی و همچنین بعد آنها curse of dimensionality دامنگیر ما نمی‌شود (علاوه بر این دانستن ارتباط این فیچرها نیاز به علم پزشکی دارد)، پس از feature selection می‌گذریم.

نمودارهای وابستگی ca، cp و exang در [10] In تا [12] In آمده.



در این شکل دیده می‌شود که exang داشتن، احتمال بیماری قلبی را کاهش می‌دهد و برعکس.



این شکل هم نشان دهنده این است که ca خیلی کم و خیلی زیاد، احتمال ابتلا به بیماری قلبی را افزایش می‌دهد.



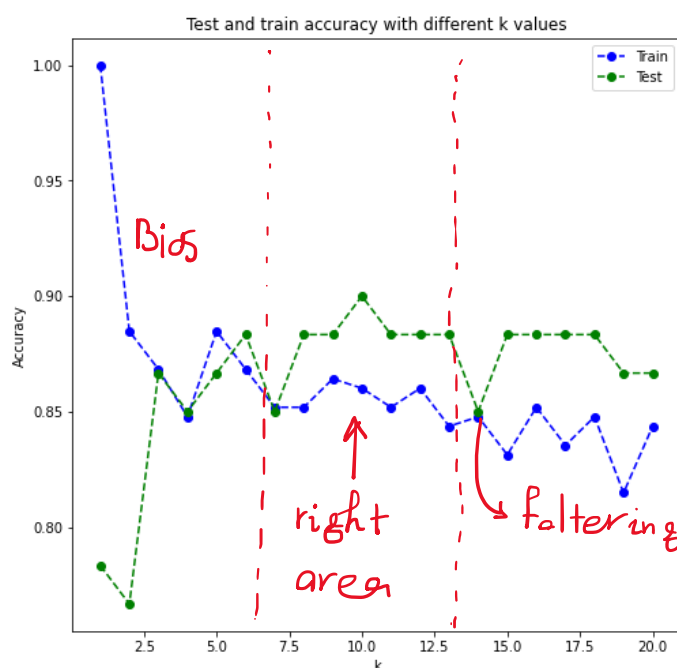
## پیش پردازش

همانطور که در [15] Out دیده می‌شود، در دیتاست داده NaN وجود ندارد، پس به مرحله‌ی بعدی می‌رویم. در [16] In برای پیش پردازش، از آنجا که scale در نتیجه kNN موثر است، نرمال‌سازی صورت گرفته. این کار به کمک StandardScaler انجام شده. نتیجه‌ی نرمال‌سازی در [16] Out قابل مشاهده است. قبل از شروع آموزش تقسیم داده‌ها صورت می‌گیرد، با شمارش نسبی داده‌های دو دسته در [19] In و [20] از تابع شافل اطمینان پیدا می‌کنیم.

## مرحله دوم | ساخت و آموزش

### kNN

در روش kNN با استفاده از  $k = 12$  به دقت 0.86 برای داده‌های آموزشی و 0.88 برای داده‌های تست رسیدیم. در [22] Out تنظیمات kNN خروجی داده شده است. نمودار دقت برای داده‌های آموزشی و تست، با  $k$  ها مختلف به شکل زیر می‌باشد:



به نظر  $k$  بین 6 و 12 برای این مسئله مناسب می‌باشد.

### Naïve Bayes

مانند kNN با استفاده از کتابخانه sklearn و با تنظیمات اولیه [27] Out به نتیجه دقت برای داده‌های آموزشی 0.83 و 0.90 برای داده‌های تست می‌رسم.

جدول مقایسه این دو روش در [30] Out به شکل زیر می باشد:

	Train	Test
kNN	86%	88%
Naïve Bayes	83%	90%

## مرحله اول | پیاده‌سازی

### Input:

$D = \{t_1, t_2, \dots, t_n\}$  // Set of elements  
 $K$  // Number of desired clusters

### Output:

$K$  // Set of clusters

### K-Means algorithm:

- 1 Assign initial values for  $m_1, m_2, \dots, m_k$
- repeat
  - 2 assign each item  $t_i$  to the clusters which has the closest mean;
  - 3 calculate new mean for each cluster;
- until convergence criteria is met;

برای الگوریتم بالا، سه تابع تعریف شده است. در [2] In تابع pick\_centroids برای مرحله‌ی یک است که با گرفتن کل نقاط و مقدار  $k$ ، مرکز اولیه بصورت رندوم از نقاط انتخاب می‌کند و پس می‌دهد.

در [3] In مرحله‌ی دو انجام شده که نزدیکترین مرکز به هر نقطه را حساب کرده و در یک numpy array شامل اندیس مرکز مرتبط بر می‌گرداند.

در [4] In مراکز آپدیت می‌شوند که مربوط به مرحله‌ی سه است. این کار با میانگین گرفتن از داده‌ها در خوشه‌ی مربوط انجام می‌شود.

الگوریتم kmeans در [5] In پیاده‌سازی شده است. این تابع داده‌ها و مقدار  $k$  و تعداد iteration ها تا پایان را از ورودی گرفته و شبه کد بالا را اجرا می‌کند

## مرحله دوم | اجرای روی دیتاست

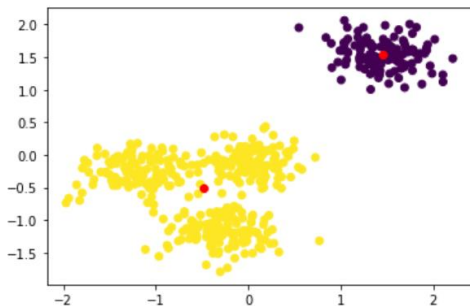
الف | در [7] In دو دیتاست داده شده، لود شده اند و scatter plot آن‌ها در [7] Out قابل مشاهده است.

\*از آنجایی که در این تمرین با داده دو بعدی کار می‌کنیم، استفاده از pandas در محاسبه فاصله‌ها سربار خیلی زیادی دارد، پس بهتر است از numpy استفاده کنیم. علاوه بر این تصویر سوال بعد هم بصورت numpy لود می‌شود.

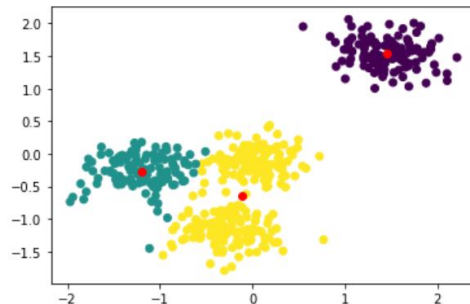
ب و ج | برای  $k = 2$  در [9] In تابع kmeans نتایج نهایی را خروجی می‌دهد (کلاستر ها و مراکزشان). در [10] Out خروجی scatter plot رسم شده و در [11] Out نتایج خطای هر کلاستر و میانگین خطای کل خوشه‌بندی چاپ شده است.

این کار برای  $k = 3$  و  $k = 4$  هم انجام شده‌اند.

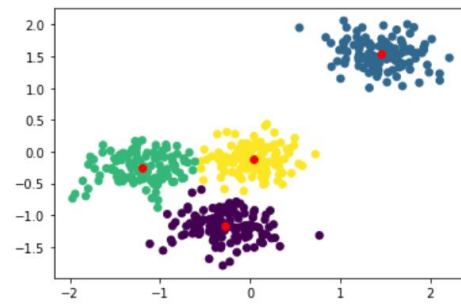
نتایج برای شروع ابتدایی ایده‌آل در شکل زیر قابل مشاهده هستند



Cluster 1 error: 0.31  
Cluster 2 error: 0.74  
Clustering Error: 0.53

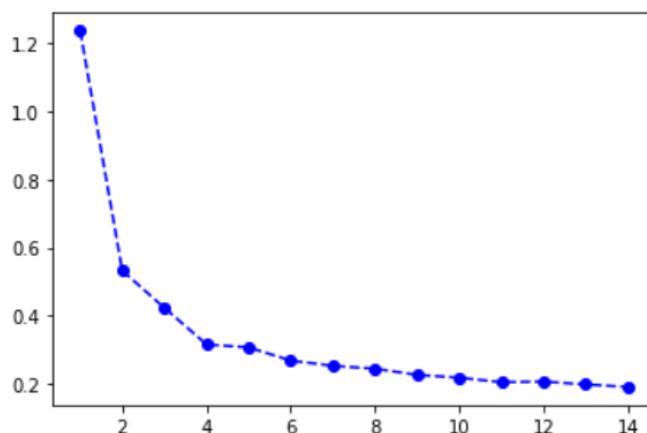


Cluster 1 error: 0.31  
Cluster 2 error: 0.33  
Cluster 3 error: 0.61  
Clustering Error: 0.42



Cluster 1 error: 0.33  
Cluster 2 error: 0.31  
Cluster 3 error: 0.32  
Cluster 4 error: 0.28  
Clustering Error: 0.31

دو هـ | در [18] In با مقادیر مختلف  $k$ ، خطای خوشه‌بندی بررسی شده و نمودار مربوطه در [18] Out قابل مشاهده است.

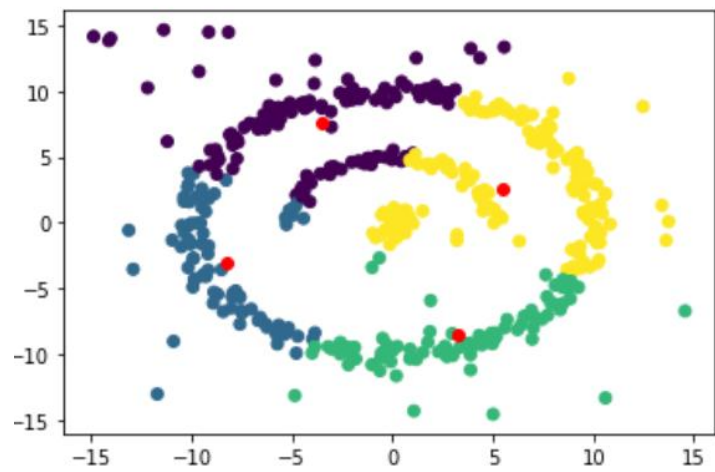


همانطور که دیده می‌شود بعد از  $k = 4$  تقریباً بهبودی نداشتیم و به elbow رسیده‌ایم. پس در این صورت  $k = 4$  برای انتخاب ایده‌آلی است.

ی | برای dataset2.csv با kmeans به نتیجه‌ی [19] Out میرسیم که قابل قبول نیست. خطای خوشه‌بندی برای این دیتاست بیشتر از 8 است ( $k = 4$ ). البته این خطا با مرحله‌ی قبلی قابل مقایسه نیست.

دلیل این اتفاق این است که kmeans بر اساس فاصله خوشه‌بندی می‌کند و خوشه‌هایش شکل دایره‌ای دارند. در dataset2 خوشه‌ها به شکل پیوسته هستند و می‌توانند با روش‌های density based یا nearest neighbor به جواب برسند.

لازم به ذکر است الگوریتم kmeans به طور کلی برای خوشه‌های با تراکم متفاوت، سایز متفاوت و شکل غیر دایره‌ای به مشکل بر می‌خورد.



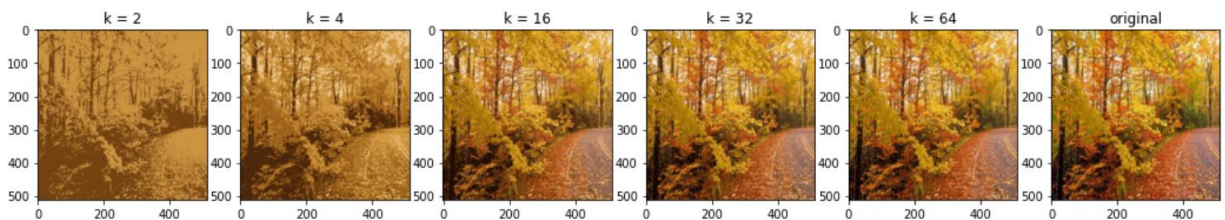
## کاهش رنگ عکس [kmeans/jupyter-image-compression.ipynb -> part 3]

در [24] In برای هر یک از اعضای لیست colors، kmeans را انجام می‌دهیم. مراحل کلی این کار برای یک k خاص به شکل زیر است:

```
1 read image as img
2 reshape img into 2d array
3 run kmeans on reshaped image: result of kmeans is a list of k colors and
  index of what each pixel belongs to
4 replace each pixel with its' center (its' representing color)
5 reshape the result into original size
```

دو بعدی کردن تصویر به این شکل انجام شده که بجای طول و عرض پیکسل، از یک ایندکس خطی استفاده شده است؛ چرا که برای کاهش رنگ نیازی به مکان پیکسل‌های نزدیک بهم نداریم و "مکان" و نزدیکی در اینجا نزدیکی رنگ‌هاست. این درحالی است که در مسائلی مثل image segmentation مکان هر پیکسل اهمیت زیادی دارد.

در نهایت خروجی تصاویر در [25] Out نمایش داده شده است.



## شناسایی نقاط پرتراکم کرونایی [dbscan/jupyter-covid.ipynb]

### مرحله اول | نمایش داده‌های خام

با کمک دستورات ذکر شده در صورت تمرین، نقاط ابتدایی در [1] In تنظیم شده اند تا نقشه به ایران اشاره کند. سپس با خواندن covid.csv، نقاط در لیست dataset ذخیره شده‌اند. در نهایت برای نمایش، هر داده با یک نقطه‌ی قرمز در آن مکان نمایش داده شده است.

### مرحله دوم | پیاده‌سازی DBSCAN

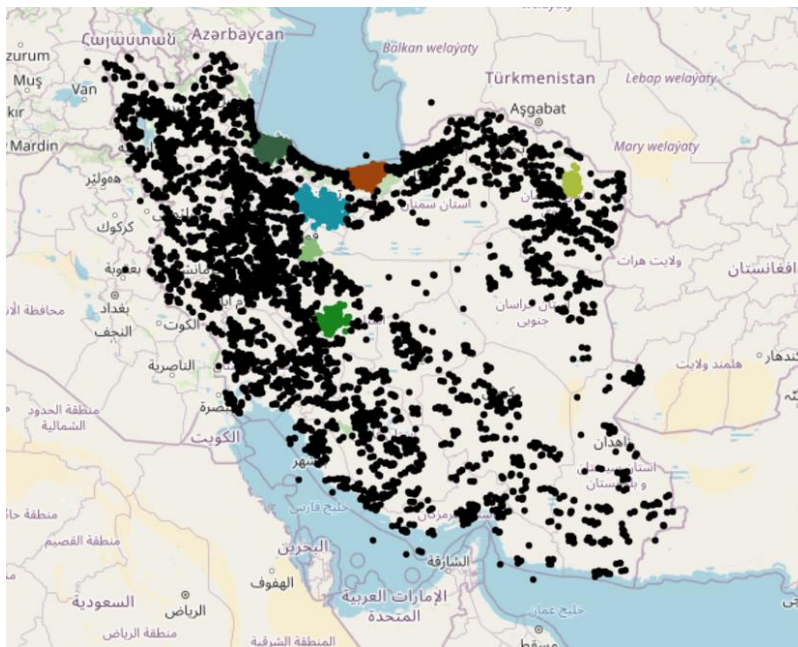
در [2] In الگوریتم DBSCAN با پارامترهای  $\text{eps} = 0.3$  و  $\text{minPts} = 10$  اجرا شده است و به 38 کلاستر و 341 داده پرت رسیدیم.

labels شامل برچسب مجازی (برای مشخص کردن داده‌های هم کلاستر) داده هاست. تعداد داده‌های پرت، تعداد داده‌هایی است که label برابر 1- دارند و سایز label های یونیک منهای یک (داده‌های پرت) تعداد کلاسترهای پیدا شده را نمایش می‌دهد.

### مرحله سوم | پیدا کردن شهرهای پرتراکم

برای پیدا کردن ۶ شهر پرتراکم، با تغییر  $\text{eps}$  و  $\text{minPts}$  به مقدار ۰.۳ برای  $\text{eps}$  و ۳۳۰ برای  $\text{minPts}$  رسیدیم. ۶ شهر پرتراکم شامل تهران، قم، مشهد، مازندران، رشت و یک شهر در جنوب قم (فکر می‌کنم کرمان باشد) است.

در نهایت خروجی نقشه در [4] Out نمایش داده شده است.



## مرحله اول | دانلود و خواندن فایل داده‌ها

این مرحله با wget! و استفاده از کتابخانه zipfile انجام شد و فایل‌ها ذخیره شدند. پس از دانلود و unpack کردن، باید داده‌ها را وارد پایتون بکنیم. برای این کار در [3] In، تک تک خطوط را خوانده و کلمات و بردار-های بازنمایی متناظر را در دو لیست ذخیره می‌کنیم. لیست شامل کلمات embedding\_words و لیست شامل بردارهای بازنمایی embedding\_vectors نام دارند. برای سادگی لیست دوم را به آرایه numpy تبدیل می‌کنیم.

## مرحله دوم | پیاده‌سازی PCA

تابع pca در [4] In تعریف شده در این مرحله با ورودی embedding\_vectors و تعداد pc های موردنیاز، مهم‌ترین بردارهای اساسی را بازمیگرداند. برای این کار از فرایند زیر پیروی می‌شود:

```
1 normalize data
2 calculate covariance matrix
3 extract biggest r eigen vectors
4 map data onto biggest r eigen vectors
```

در نهایت با اجرای تابع بالا با ورودی embedding\_vectors و  $r=2$  خروجی مورد نظر، یعنی داده‌ها با ابعاد کمتر بدست می‌آیند.

## مرحله سوم | نمایش داده‌های کاهش بعد داده شده

در این مرحله در [7] In با انتخاب کلمات قابل تمایز، به نتیجه‌ی زیر رسیدیم:

