

i	$x_1^{(i)}$	$x_2^{(i)}$		
1	0	-1		
2	1	0		
3	2	1	μ_1	
4	1	1	$x_1 - \mu_1$	
5	-1	1	$x_2 - \mu_2$	
6	-1	-1		
7	-1	-1		

$$= \tilde{X} = \begin{bmatrix} -1/7 & -1 \\ 1/7 & 0 \\ 13/7 & 1 \\ 6/7 & 1 \\ -8/7 & 1 \\ -8/7 & -1 \\ -8/7 & -1 \end{bmatrix}$$

$$\mu_1 = \frac{0+1+2+1+1-1-1}{7} = \frac{1}{7}$$

$$\mu_2 = \frac{-1+0+1+1+1-1-1}{7} = 0$$

$$\Rightarrow \mu = \begin{bmatrix} 1/7 \\ 0 \end{bmatrix}$$

$$C = \frac{1}{7} \tilde{X}^T \tilde{X} = \frac{1}{7} \begin{bmatrix} \sum_{i=1}^7 \tilde{x}_1^{(i)2} & \sum_{i=1}^7 \tilde{x}_1^{(i)} \tilde{x}_2^{(i)} \\ \sum_{i=1}^7 \tilde{x}_1^{(i)} \tilde{x}_2^{(i)} & \sum_{i=1}^7 \tilde{x}_2^{(i)2} \end{bmatrix}$$

$$= \frac{1}{7} \begin{bmatrix} \frac{1+36+169+36+14+64+64}{49} & \frac{1+0+13+6-8+8+8}{7} \\ \frac{1+0+13+6-8+8+8}{7} & 1+0+1+1+1+1+1 \end{bmatrix}$$

$$= \frac{1}{7} \begin{bmatrix} \frac{62}{7} & 4 \\ 4 & 6 \end{bmatrix} = \begin{bmatrix} \frac{62}{49} & \frac{4}{7} \\ \frac{4}{7} & \frac{6}{7} \end{bmatrix}$$

$$\det(C - \lambda I) = 0 \rightarrow \begin{vmatrix} \frac{62}{49} - \lambda & \frac{4}{7} \\ \frac{4}{7} & \frac{6}{7} - \lambda \end{vmatrix} = 0$$

$$\lambda^2 - \frac{104}{49} \lambda + \frac{260}{343} = 0$$

$$\rightarrow \begin{cases} \lambda_1 = \frac{52 + 2\sqrt{221}}{49} \approx 1.67 \\ \lambda_2 = \frac{52 - 2\sqrt{221}}{49} \approx 0.45 \end{cases}$$

$$CV_1 = \lambda_1 V_1$$

$$\begin{bmatrix} \frac{62}{49} v_1 + \frac{4}{7} v_2 \\ \frac{4}{7} v_1 + \frac{6}{7} v_2 \end{bmatrix} = \begin{bmatrix} 1.67 v_1 \\ 1.67 v_2 \end{bmatrix}$$

↓

$$\approx \begin{bmatrix} 1.26 v_1 + 0.57 v_2 \\ 0.57 v_1 + 0.86 v_2 \end{bmatrix} \Rightarrow \begin{aligned} -0.41 v_1 &= -0.57 v_2 \\ 0.57 v_1 &= +0.81 v_2 \end{aligned}$$

$$\Rightarrow v_2 = 0.71 v_1 \Rightarrow \begin{bmatrix} 1.4 \\ 1 \end{bmatrix} = V_1 \rightsquigarrow \text{1st PC}$$

$$CV_2 = \lambda_2 V_2$$

$$\begin{bmatrix} 1.26 v_1 + 0.57 v_2 \\ 0.57 v_1 + 0.86 v_2 \end{bmatrix} = \begin{bmatrix} 0.45 v_1 \\ 0.45 v_2 \end{bmatrix} \Rightarrow \begin{aligned} 1.81 v_1 &= -0.57 v_2 \\ 0.57 v_1 &= -0.41 v_2 \end{aligned}$$

$$\Rightarrow v_2 = -1.9 v_1 \Rightarrow \begin{bmatrix} -0.7 \\ 1 \end{bmatrix} = V \rightsquigarrow \text{2nd PC}$$

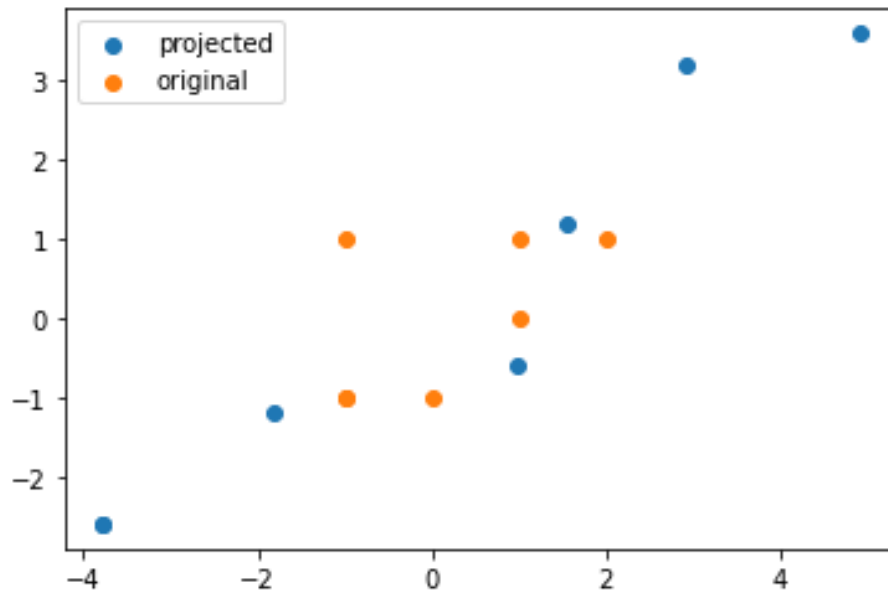
$$A = V_1$$

$$X' = \tilde{X}A = \begin{bmatrix} -\frac{1}{7} \times \frac{14}{10} - 1 & = & -1.2 \\ \frac{6}{7} \times \frac{14}{10} - 0 & = & 1.2 \\ \frac{13}{7} \times \frac{14}{10} + 1 & = & 3.6 \\ \frac{6}{7} \times \frac{14}{10} + 1 & = & 3.2 \\ -\frac{8}{7} \times \frac{14}{10} + 1 & = & -0.6 \\ -\frac{8}{7} \times \frac{14}{10} - 1 & = & -2.6 \\ -\frac{2}{7} \times \frac{14}{10} - 1 & = & -2.6 \end{bmatrix}$$

$$A^T = V_1^T = \begin{bmatrix} 1.4 & 1 \end{bmatrix}$$

$$X_{re} = X A^T - \begin{bmatrix} \mu_1 & \mu_2 \\ \mu_1 & \mu_2 \\ \mu_1 & \mu_2 \\ \mu_1 & \mu_2 \\ \mu_1 & \mu_2 \\ \mu_1 & \mu_2 \\ \mu_1 & \mu_2 \end{bmatrix}_{7 \times 2}$$

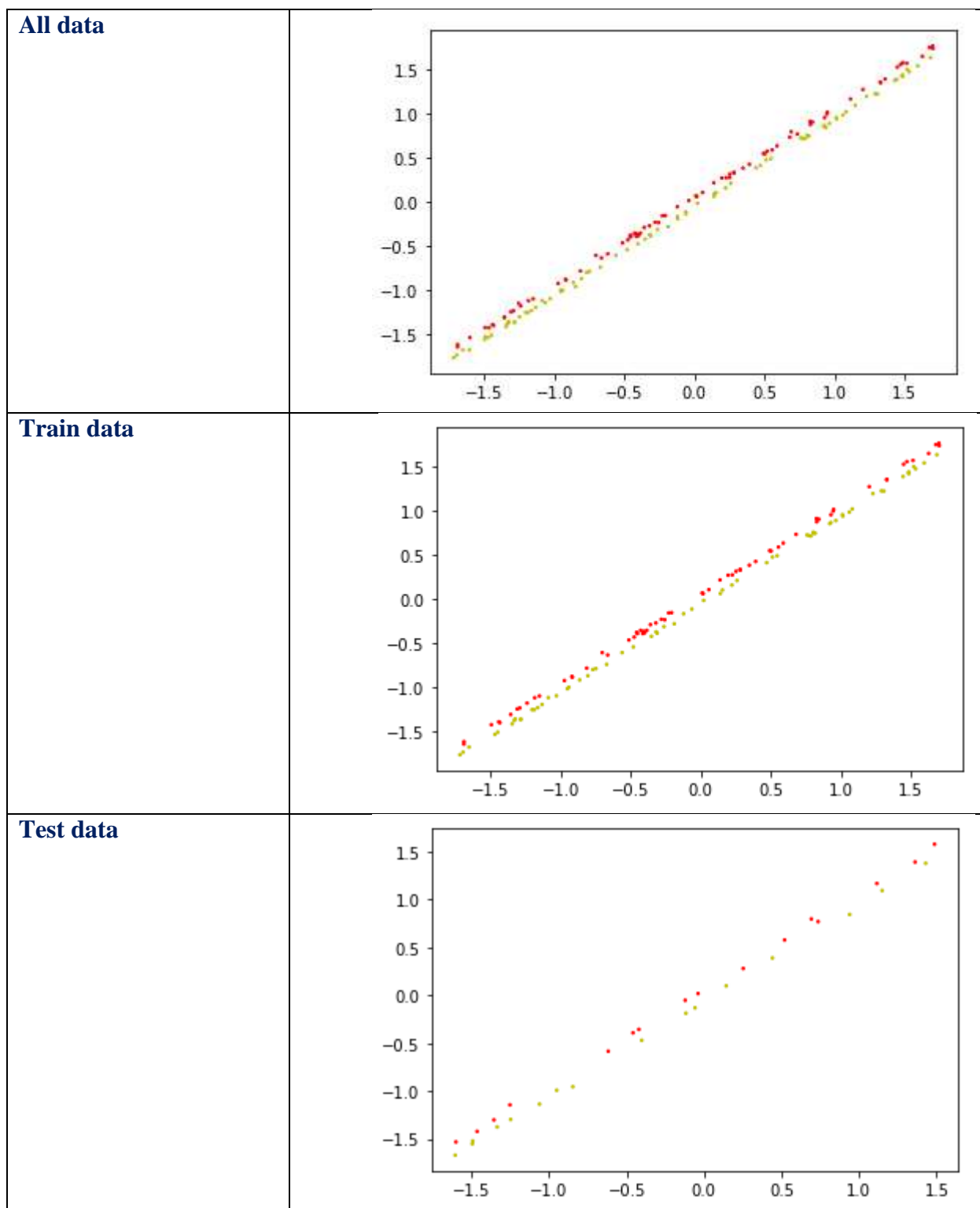
$$= \begin{bmatrix} -1.68 - \mu_1 & -1.2 - \mu_2 \\ 1.68 - \mu_1 & 1.2 - \mu_2 \\ 5.04 - \mu_1 & 3.6 - \mu_2 \\ 3.08 - \mu_1 & 3.2 - \mu_2 \\ -0.84 - \mu_1 & -0.6 - \mu_2 \\ -3.64 - \mu_1 & -2.6 - \mu_2 \\ -3.64 - \mu_1 & -2.6 - \mu_2 \end{bmatrix} = \begin{bmatrix} -1.82 & -1.2 \\ 1.54 & 1.2 \\ 4.9 & 3.6 \\ 2.93 & 3.2 \\ 0.98 & -0.6 \\ -3.78 & -2.6 \\ -3.78 & -2.6 \end{bmatrix}$$



سوال ۲:

ویژگی ها و برچسب های هردو دسته دادگان تست و آموزش را به صورت ماتریس های داده ی `numpy` خواندم. برچسب های ۰ و ۱ را به ۱- و ۱ تبدیل کردم (فقط برای راحتی خطاب داده ها) و کلاس ۱- کلاس قرمز و کلاس ۱+ را کلاس زرد می نامم. قبل از انجام هرکار داده ها را با `StandardScaler` از `sklearn` استاندارد می کنم. این استانداردسازی با برآورد میانگین و واریانس داده های آموزش به عنوان میانگین و واریانس مجهول توزیع انجام شده. و روی داده های تست هم همین مدل استفاده می شود.

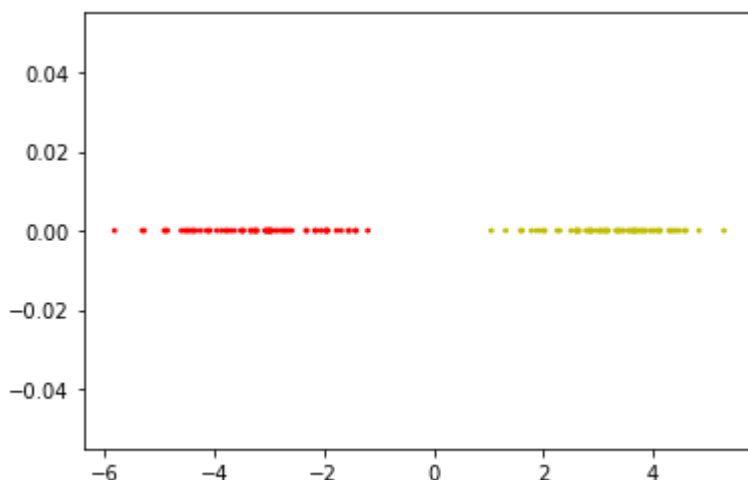
در قدم اول نمایش داده های تست و آموزش و هر دو باهم :



با استفاده از تابع LinearDiscriminantAnalysis از مجموعه کتابخانه ی sklearn یک مدل LDA روی داده های آموزش
برازش می‌دهم. بردار راستای LDA به دست می‌آید :

[386.49722, -386.96594]

که داده های دوبعدی را به فضای تک بعد می برد. داده های آموزش روی این فضای تک بعدی:



واضح است که داد ها را می توان کاملا با تعیین یک عدد آستانه از هم تفکیک کرد. برای تعیین این مقدار آستانه از یک پرسپترون باز هم از کتابخانه ی sklearn استفاده می کنم. این مدل perc1 بعد از برازش روی داده های آموزش، برابر با بردار زیر می شود:

[1. 2.49924668]

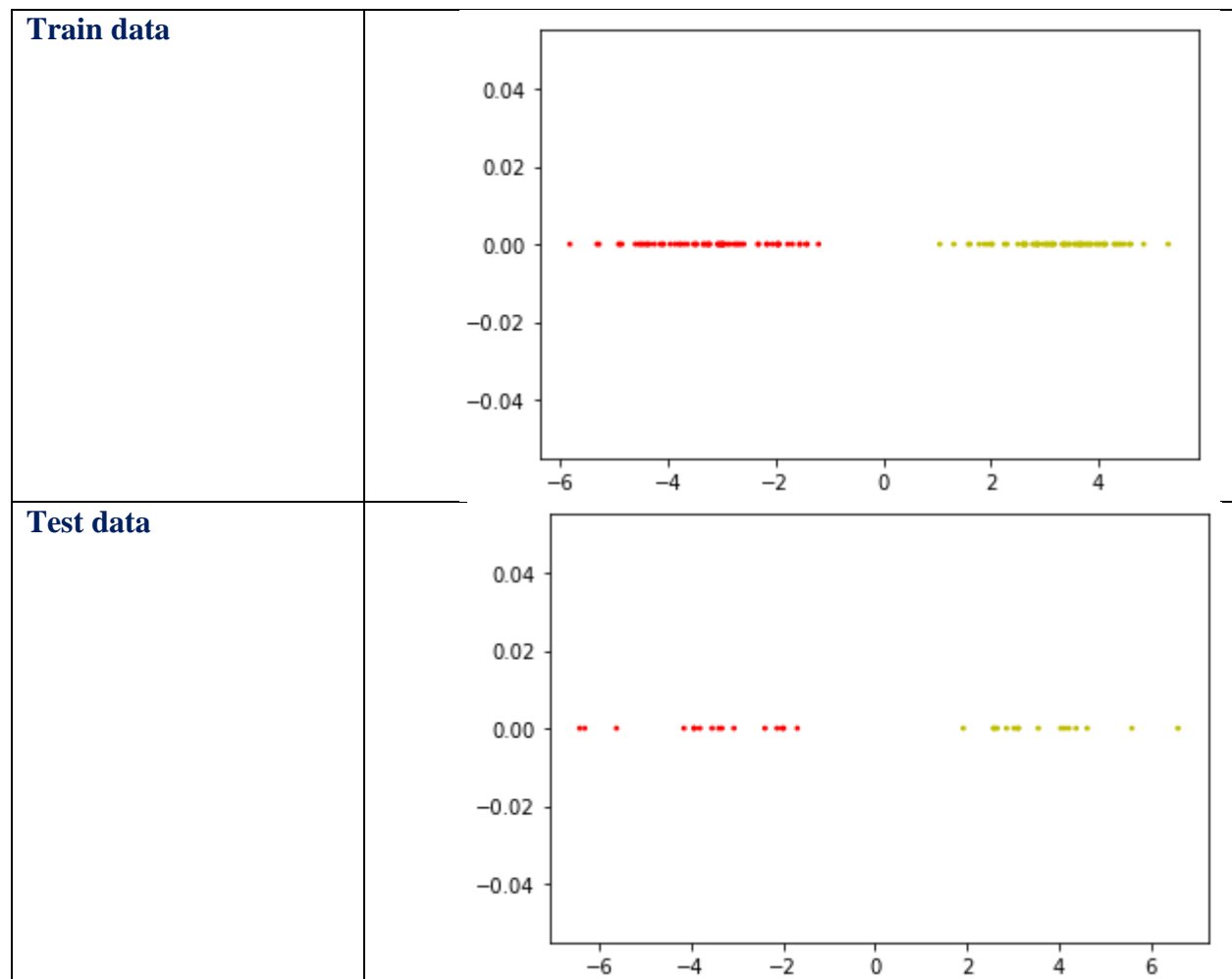
عدد ۱ بایاس مدل است. حالا با برازش این بردار مدل روی داده های آموزش دقت آموزش برابر می شود با :

	precision	recall	f1-score	support
-1	1.00	1.00	1.00	64
1	1.00	1.00	1.00	64
accuracy			1.00	128
macro avg	1.00	1.00	1.00	128
weighted avg	1.00	1.00	1.00	128

و بعد از تبدیل داد های تست با بردار راستای LDA و دسته بندی نتایج با perc1 روی داده های تست دارم :

	precision	recall	f1-score	support
-1	1.00	1.00	1.00	16
1	1.00	1.00	1.00	16
accuracy			1.00	32
macro avg	1.00	1.00	1.00	32
weighted avg	1.00	1.00	1.00	32

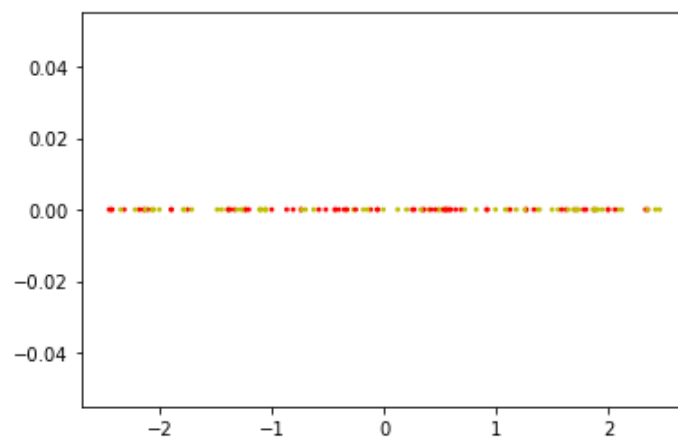
دقت دسته بند دو کلاسه روی هر دو دادگان تست و آموزش برابر ۱۰۰ درصد هست و خطای مدل روی این داده ها صفر است. اگر داده های کلاس بندی شده ی یک بعدی را روی شکل نشان دهم ، هیچ داده ای که خطا کلاس بندی شده باشد (رنگ سیاه) ندارم. دو کلاس فاصله ی زیاد از هم دارند.



با استفاده از مدل PCA از کتابخانه ی sklearn داده های آموزش را با معیار PCA به یک فضای تک بعدی منتقل می کنم.
ماتریس کوواریانس این تبدیل :

```
[ [1.0078739  1.0062541] ,  
  [1.0062541  1.0078735] ]
```

این داده های به دست آمده را روی محور بر اساس رنگ کلاس هستند :



این تبدیل داده های کلاس را از هم جدا نمی کند بلکه بیشتر از قبل آن ها را توی هم ادغام می کند. بازهم با یک مدل پرسپترون داده های آموزش را مدل می کنیم. مدل perc2 به دست آمده :

[0. -0.5916291]

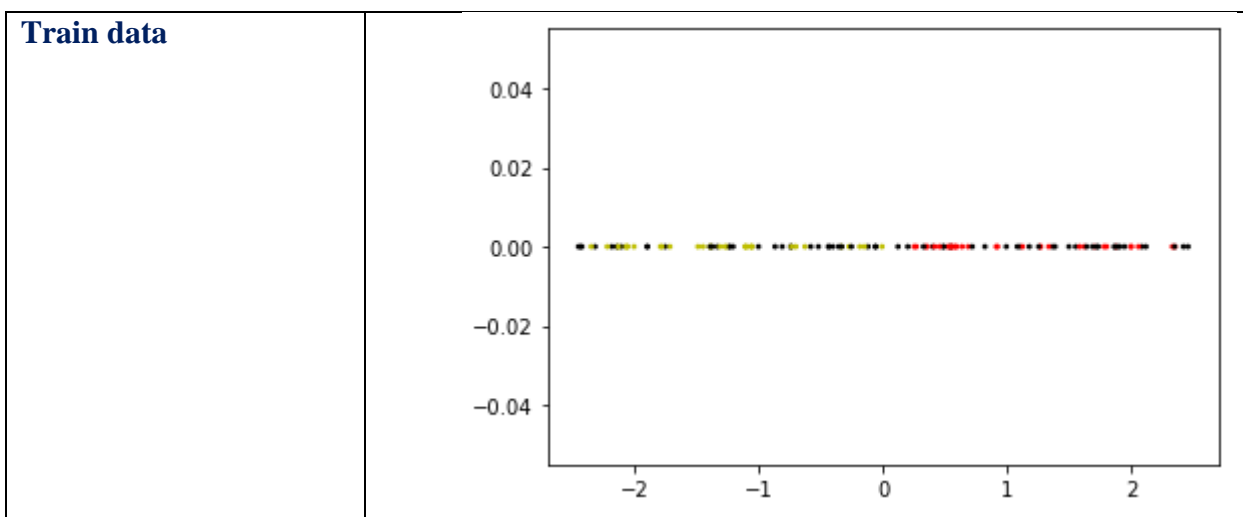
داده های آموزش تک بعدی شده با مولفه ی PCA را با این مدل کلاس بندی میکنم. نتایج زیر به دست می آید.

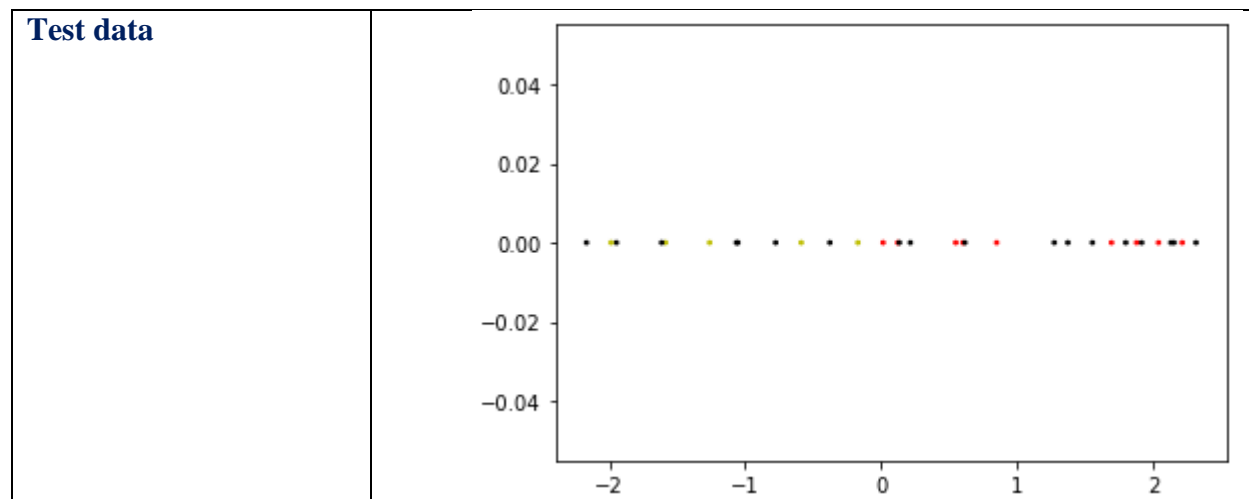
	precision	recall	f1-score	support
-1	0.50	0.50	0.50	64
1	0.50	0.50	0.50	64
accuracy			0.50	128
macro avg	0.50	0.50	0.50	128
weighted avg	0.50	0.50	0.50	128

دقت مدل روی مجموعه ی آموزش پایین است و تنها نیمی از داده ها در کلاس درست طبقه بندی می شوند. همچنین اگر معیار PCA را روی داده های تست اعمال و سپس با همین مدل دسته بندی کنم :

	precision	recall	f1-score	support
-1	0.45	0.56	0.50	16
1	0.42	0.31	0.36	16
accuracy			0.44	32
macro avg	0.43	0.44	0.43	32
weighted avg	0.43	0.44	0.43	32

دقت مدل از همان ۵۰ درصد هم پایین تر می آید. همچنین روی شکل :





اگر اولین نمایش داده ها در دو بعد را دوباره نگاه کنیم، روی همان نمودار جداسازی بهتری بین داده های می شد متصور بود تا روی داده های انتقال یافته با PCA. این مدل فاقد کارایی است.

تبدیل PCA برای انتقال داده ها و کاهش بعد با هدف بیشینه سازی واریانس داده های باقیمانده پیش می رود و تلاش می کند داده ها را در جهتی نگاشت دهد که حداکثر واریانس در آن وجود دارد و از طرفی ابعادی که تقریباً نمودار یکنواخت و مسطحی دارند را حذف کند. PCA به برچسب داده ها کاملاً بی توجه است (روش بدون نظارت) در حالی که معیار LDA به طور خاص روی داده های برچسب دار با هدف تشدید تمایز بین دو کلاسی در بعد پایین تر انتقال را انجام می دهد. از روی همین می توان انتظار داشت که اگر هدف دسته بندی داده ها باشد، معیار LDA نتیجه ی مطلوب تری خواهد داشت. این نتیجه روی مدل هایی که بالا ساخته شد به خوبی مشخص شد. اگر به اولین نمایش داده ها (دو بعدی) برگردیم، این دو کلاس روی همان دو بعد از هم تفکیک پذیرند اما فاصله ی دو کلاس از خط جداکننده ی فرضی کوچک است. کاری که LDA انجام داد تشدید همین فاصله بود طوری که روی محور اعدادی که LDA ساخت، فاصله ی آخرین داده های دو کلاس حدود ۲ است در حالی که تمرکز داده های هر کلاس را هم بالاتر برد که این تمرکز روی نمایش دوبعدی تقریباً اصلاً دیده نمی شد. معیار LDA واریانس داده ها را به طور انتخابی و با توجه به برچسب کلاسه ها کم یا زیاد می کند در حالی که PCA در جهت حفظ واریانس کل داده ها تلاش می کند که روی این داده های نزدیک به هم (ولی جدایی پذیر) باعث ادغام کامل دسته ها شد چون بیشترین واریانس داده های اولیه نه در جهتی که آن دو کلاس از هم متمایز می شد (خط فرضی بین دو کلاس) که در همان جهتی تقریباً عمود بر آن بود. یعنی واریانس داده ها در جهت اشتباهی برای هدف جداسازی حفظ شد چون PCA بردار مولفه اول را انتخاب می کند در حالی که آنچه در این دسته بندی نیاز است، حفظ واریانس در جهت بردار مولفه ی دوم است و به همین دلیل است که در PCA کلاس ها جداناپذیر شد.

سوال ۳:

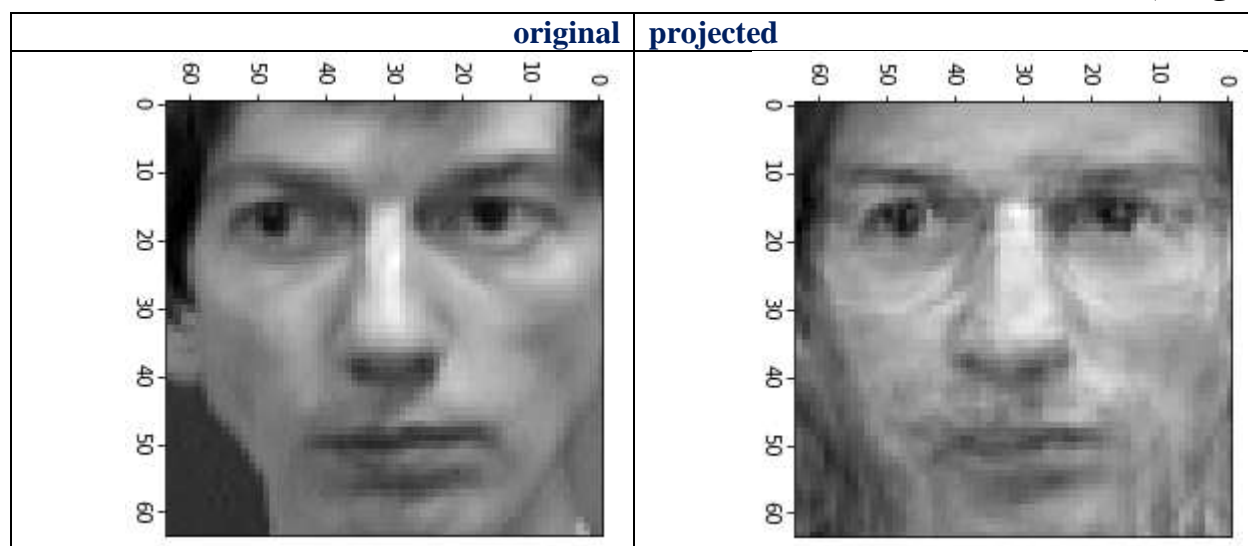
فایل را با `scipy` خواندم. برای اینکه راحت تر با داده ها کار کنم بخش داده هایی که تیار دارم یعنی `faces` را که مجموعاً ۴۰۰ عکس و هر عکس دارای ۶۴*۶۴ مولفه بود را به صورت یک ماتریس داده ی `numpy` از ۴۰۰ ردیف و ۶۴*۶۴ ستون ذخیره کردم. این عکس ها به ترتیب ۱۰ تا ۱۰ عکس های ۴۰ نفر اند پی یک بردار ۴۰۰ تایی هم برای برچسب ها درست کردم که البته چون کار دسته بندی و شناسایی چهره خواسته نشده بود به کار نیامد! این ۴۰۰ ردیف را بدون `shuffle` کردن با `train_test_split` از `sklearn` به دو دسته ۲۸۰ تایی و ۱۲۰ به ترتیب به عنوان آموزش و تست تقسیم کردم.

با یک مدل PCA از sklearn ۵ این داده ها را از فضای ۶۴*۶۴ به ۴۵ بعد می بریم. مدل فقط با ۲۸۰ داده ی اول آموزش دیده. سپس داده های آموزش و تست را با مولفه های به دست آمده به فضای ۴۵ بعد می بریم و بعد دوباره به همان فضای اولیه بازسازی می کنیم. مقدار RMSE به دست آمده بین داده های اولیه ی هر مجموعه با فضای بازسازی شده برابر است با :

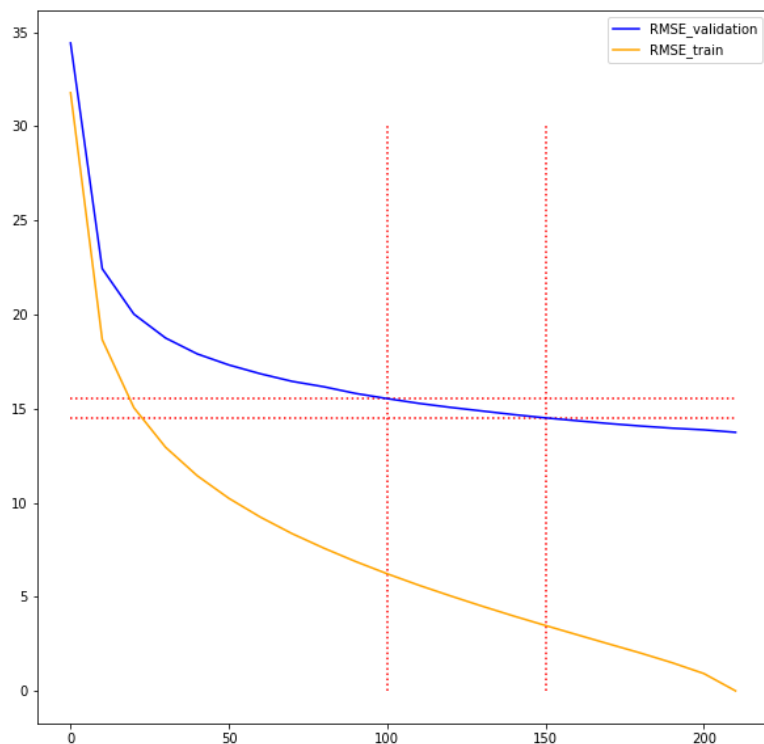
RMSE on train data = 11.525403972573551

RMSE on test data = 18.023206734576966

در زیر یک تصویر دلخواه از مجموعه ی تست را (تصویر نظیر ردیف 1 تست) در دو حالت اولیه و بازسازی شده با ۴۵ مولفه ی اصلی داریم.

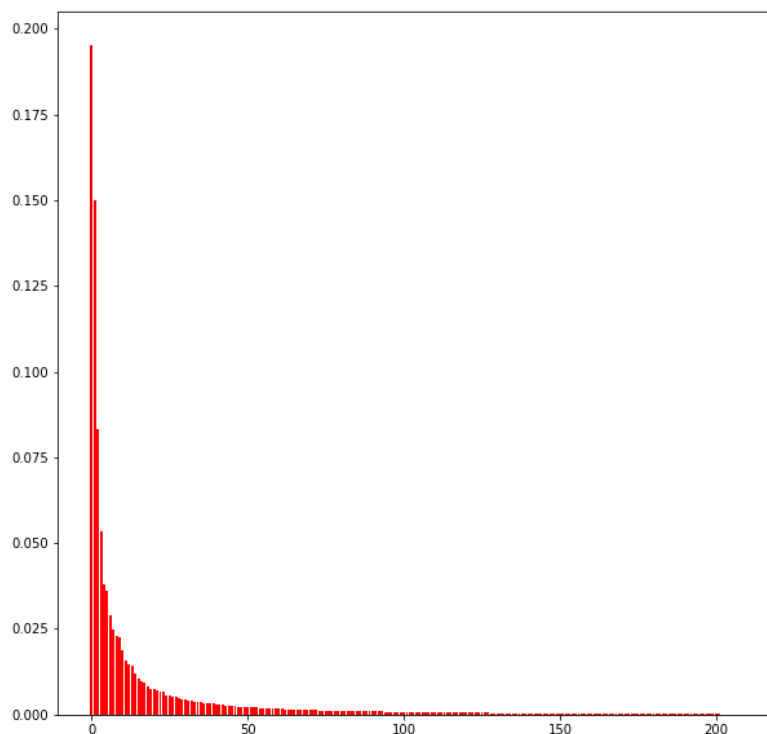


برای انتخاب بعد بهینه مجموعه ی آموزش را به اندازه ی ۷۵ به ۲۵ تقسیم می کنیم. و بخش آخر را برای validation نگه می دارم. برای فضای ثانویه از کل مجموعه ی ۲۱۰ تایی تست یک مجموعه ی بزرگ از ۱ تا ۲۱۰ می توانم به عنوان عدد بعد انتخاب کنم (اگر داده ها بیشتر بود، می توانستم تا ۶۴*۶۴ را به عنوان بعد ثانویه بگیرم ولی اینجا محدودیت تعداد برای ماتریس داده ها دارم) با شروع از صفر و ۱۰ تا ۱۰ جلو رفتن هر بار یک تبدیل PCA روی مجموعه ی تست می سازم. به ازای هر کدام از این تبدیل های ساخته شده روی مجموعه های ۲۱۰ داده ی آموزش و ۷۰ داده ی اعتبارسنجی مقدار تبدیل را محاسبه و سپس فضای اولیه را با همان بردارها بازسازی می کنیم. RMSE را محاسبه می کنیم. نمودار این اعداد به دست آمده :

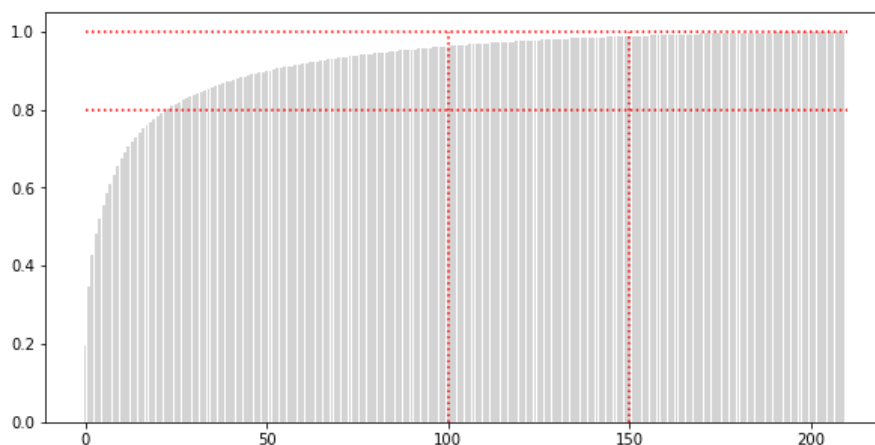


اگر روی بخش مشخص شده ی نمودار با خط چین های قرمز دقت کنم اولاً از ۱۰۰ به بعد شیب خط آبی اعتبارسنجی ملایم تر می شود و بهبود بیشتر روی داده های تست است. پس عدد بعد را می توان بعد از ۱۰۰ انتخاب کرد. از طرفی بعد از ۱۵۰ هم شیب زیادی ملایم می شود. و ۵۰ بعد باقیمانده بیشتر فقط حدود ۱ واحد روی RMSE تاثیر دارد. از این مقدار می شود چشم پوشید. پس عدد را جایی بین ۱۰۰ تا ۱۵۰ انتخاب می کنم.

اگر روی نمودار میله ای نسبت واریانس های توزیع یافته به ازای هر بردار مولفه ی ۱ تا ۲۱۰ دقت کنیم کاهش نسبت از ۱۰۰ به بعد به سختی قابل مشاهده است.



برای دید بهتر نمودار همین نسبت ها را به صورت تجمعی رسم می‌کنم. از بردار ۱۵۰ به بعد تاثیر بسیار ناچیز است. و آخرین بردارهایی که تاثیر ثابل توجه دارند همان ۱۰۰ تا ۱۵۰ هستند. روی همان بازه ای که در نظر داشتیم نقطه ی وسط را به عنوان بعد ثانویه انتخاب می‌کنم. البته خود ۱۰۰ یا ۱۵۰ هم برای انتخاب مناسب است چون به لحاظ دقت بازهم تفاوت چندانی روی این بازه برای این تعداد داده نداریم. من ۱۲۵ را انتخاب می‌کنم.



انتظار دارم نتیجه ی تبدیل PCA نهایی که با بعد ۱۲۵ می‌سازم روی داده های تست مقداری حدود ۱۵ از روی نمودار اول مثل داده های اعتبارسنجی داشته باشد. آخرین مدل را با ۱۲۵ مولفه ی اصلی می‌سازم. نتایج :

```
RMSE on train data = 4.775306976555008
RMSE on validation data = 14.968495895510301
RMSE on test data = 16.10800194377165
```

دقت بازسازی روی تست هم همان نزدیک مقدار انتظار است. خطا ی بین داده های آموزش اولیه و بازسازی شده روی داده های آموزش خیلی کم است.

$$X = \begin{bmatrix} & \\ & \\ & \\ & \\ & \\ & \\ & \\ & \\ & \\ & \end{bmatrix}_{n \times d}$$

نقش: فضای کم بعدی n ردیف داده
هدف: کمابست داده ها بر فضای m بعدی با هدف
بیشینه سازی واریانس که $m < d$

اگر فرض بگیریم: $m=1$ و بردار V_1 یکد با λ_1 بزرگ است از جهت فضای ثانویه بگیریم:

$$V_1^T V_1 = 1$$

$$x^{(i)} = V_1^T x^{(i)} e^{d^1}$$

و اگر μ برابر میانگین داده ها در فضای کم بعدی اول و μ' میانگین داده ها در فضای $m=1$ بعدی دوم باشد:

$$\mu' = \frac{1}{n} \sum_{i=1}^n x^{(i)'} = \frac{1}{n} \sum_{i=1}^n V_1^T x^{(i)} = V_1^T \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} \right)$$

$$= V_1^T \mu$$

محاسبه برای واریانس داده های فضای ثانویه داریم:

$$\text{Var}(X') = \frac{1}{n} \sum_{i=1}^n (x^{(i)'} - \mu')^2 = \frac{1}{n} \sum_{i=1}^n (V_1^T x^{(i)} - V_1^T \mu)^2$$

$$= \frac{1}{n} \sum_{i=1}^n (V_1^T x^{(i)} - V_1^T \mu)(V_1^T x^{(i)} - V_1^T \mu)^T$$

$$= \frac{1}{n} \sum_{i=1}^n V_1^T (x^{(i)} - \mu) (x^{(i)} - \mu)^T V_1$$

$$= V_1^T \left(\frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T \right) V_1 = V_1^T S V_1 = J(V_1)$$

ماتریس کوواریانس فضای کم بعدی

$$\max J(V_1) = V_1^T S V_1$$

$$\text{s.t. } \|V_1\| = 1$$

حالا صحت هدف، من خواهم واریانس X' را بیشینه کنم:
تعیین کردن: این دلیل افراطی هدف بیشینه سازی
به سمت بزرگ کردن $J(V_1)$ از راه بزرگ کردن اندازه V_1 منور. من توان این قید را

تابع هزینه در رابطه J تبدیل کرد:

$$\max J(V_1) = V_1^T S V_1 + \lambda_1 (1 - V_1^T V_1)$$

$$\text{ح} \Rightarrow \frac{\partial J}{\partial V_1} = 2S V_1 - 2\lambda_1 V_1 = 0$$

در نهایت آرایشی که در این رابطه

$$S V_1 = \lambda_1 V_1$$

مقدار λ_1 این میزان وابستگی بردار ویژه S است و مقدار ویژه λ_1 آن

در نظر $\|V_1\| = 1$

$$\times V_1^T \rightarrow S V_1 = \lambda_1 V_1$$

$$V_1^T S V_1 = V_1^T (\lambda_1 V_1) = \lambda_1 (V_1^T V_1)$$

$$V_1^T S V_1 = \lambda_1 \xrightarrow{\text{Var}(X') = V_1^T S V_1} \text{Var}(X') = \lambda_1$$

در این V_1 که $\frac{\partial J}{\partial V_1} = 0$ مقدار $J(V_1) = \text{Var}(X')$ برابر می شود با مقدار ویژه ای از S که بیشترین مقدار را می تواند بگیرد برابر می شود با $\max\{\lambda_1, \lambda_2, \dots\}$ که تمام λ_i ها مقادیر ویژه S هستند و بزرگ λ_1 و λ_2 و ... هم می توان به ترتیب V_1 و V_2 را انتخاب کرد و بزرگترین λ را انتخاب کرد (معیار انتخاب PCA)