

i	$x_1^{(i)}$	$x_2^{(i)}$	$y^{(i)}$
1	0	-1	0
2	1	0	0
3	2	1	0
4	1	1	1
5	-1	1	1
6	-1	-1	1
7	-1	-1	1

$N_0 = 3$

$N_1 = 4$

$$W = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$w^T x^{(i)} = w_1 x_1^{(i)} + w_2 x_2^{(i)}$$

$$\mu_0 = \frac{1}{N_0} \sum_{y^{(i)}=0} x^{(i)} = \frac{1}{N_0} \begin{bmatrix} \sum_{y^{(i)}=0} x_1^{(i)} \\ \sum_{y^{(i)}=0} x_2^{(i)} \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 0+1+2 \\ -1+0+1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\mu_1 = \frac{1}{N_1} \sum_{y^{(i)}=1} x^{(i)} = \frac{1}{N_1} \begin{bmatrix} \sum_{y^{(i)}=1} x_1^{(i)} \\ \sum_{y^{(i)}=1} x_2^{(i)} \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 1-1-1-1 \\ 1+1-1-1 \end{bmatrix} = \begin{bmatrix} -1/2 \\ 0 \end{bmatrix}$$

$$\mu'_0 = w^T \mu_0 = w_1 \mu_{01} + w_2 \mu_{02} = w_1$$

$$\mu'_1 = w^T \mu_1 = w_1 \mu_{11} + w_2 \mu_{12} = -\frac{1}{2} w_1$$

$$S_0'^2 = \sum_{y^{(i)}=0} (w^T x^{(i)} - \mu'_0)^2 = \sum_{y^{(i)}=0} (w_1 x_1^{(i)} + w_2 x_2^{(i)} - w_1)^2 =$$

$$(0 - w_2 - w_1)^2 + (w_1 + 0 - w_1)^2 + (2w_1 + w_2 - w_1)^2 = 2(w_1 + w_2)^2$$

$$S_1'^2 = \sum_{y^{(i)}=1} (w^T x^{(i)} - \mu'_1)^2 = \sum_{y^{(i)}=1} (w_1 x_1^{(i)} + w_2 x_2^{(i)} + \frac{1}{2} w_1)^2 =$$

$$(w_1 + w_2 + \frac{1}{2} w_1)^2 + (-w_1 + w_2 + \frac{1}{2} w_1)^2 + (-w_1 - w_2 + \frac{1}{2} w_1)^2 \times 2 =$$

$$(\frac{3}{2} w_1 + w_2)^2 + (\frac{1}{2} w_1 - w_2)^2 + 2(\frac{1}{2} w_1 - w_2)^2 = 3w_1^2 + 4w_2^2 + 4w_1 w_2$$

$$J(w_1, w_2) = \frac{(\mu'_0 - \mu'_1)^2}{S_0'^2 + S_1'^2} = \frac{\frac{3}{4} w_1^2}{5w_1^2 + 6w_2^2 + 8w_1 w_2}$$

$$\frac{\partial J}{\partial w_1} = \frac{9}{4} \times \frac{2w_1(5w_1^2 + 6w_2^2 + 8w_1w_2) - w_1^2(10w_1 + 8w_2)}{(5w_1^2 + 6w_2^2 + 8w_1w_2)^2} = 0$$

$$10w_1^3 + 12w_1w_2^2 + 16w_1^2w_2 - 10w_1^3 - 8w_1^2w_2 = 0$$

$$4w_1w_2(3w_2 + 4w_1 - 2w_1) = 0$$

$$\rightarrow \begin{cases} w_1 = 0 \\ w_2 = 0 \\ 3w_2 = -2w_1 \end{cases}$$

$$\frac{\partial J}{\partial w_2} = \frac{9}{4} \times \frac{0 - w_1^2(0 + 12w_2 + 8w_1)}{(5w_1^2 + 6w_2^2 + 8w_1w_2)^2} = 0$$

$$w_1^2(12w_2 + 8w_1) = 0$$

$$\rightarrow \begin{cases} w_1 = 0 \\ 3w_2 = -2w_1 \end{cases}$$

$$3w_2 = -2w_1 \quad \leadsto \quad \frac{4}{9}w_1^2 + w_2^2 = 1 \quad \leadsto \quad w_1 = \frac{3}{\sqrt{13}}$$

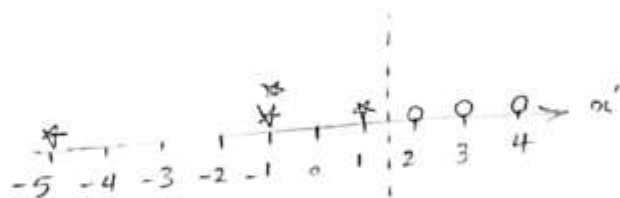
$$w_2^2 + w_1^2 = 1$$

$$w_2 = -\frac{2}{\sqrt{13}}$$

$$\uparrow \begin{bmatrix} 3 \\ -2 \end{bmatrix}$$

این بردار را به ما می‌دهد. راستی: بردار \vec{w} در $-2x = 3y$ قرار می‌گیرد. پس $x_2 = w_2$ و $x_1 = w_1$.

	x_1	x_2	$x' = wx$
1	0	-1	2
2	1	0	3
3	2	1	4
4	1	1	1
5	-1	1	-5
6	-1	-1	-1
7	-1	-1	-1



$x' \leq 1.5$ class *
 $y = 1$

$x' > 1.5$ class 0
 $y = 0$

$$f(a_1, a_2) = 3a_1 - 2a_2 + 1.5$$

مدل کلاس بندی خطی

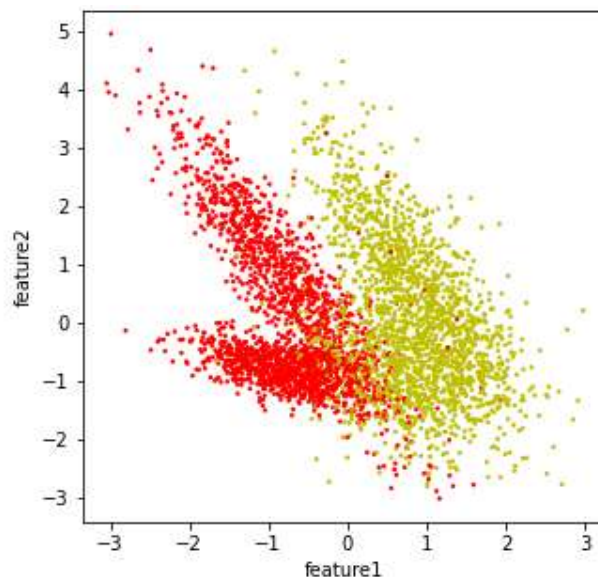
$$\frac{\text{تعداد اشتباه کلاس بندی شده}}{\text{تعداد داده ها}} = 0 \Rightarrow 100\%$$

عملکرد مدل روی داده های آموزش = 100%

تمرین دوم :

روی داده های فایل `train.csv` مدل هر مرحله ساخته شده و با داده های `test.csv` پارامترهای مدل ارزیابی شده. این مجموعه داده ها سه ستون دارد که دو ستون اول `feature1` و `feature2` به عنوان ویژگی های ورودی استفاده شده و `label` را برچسب هر داده گرفتیم که مقداری برابر ۰ یا ۱ دارد که برای راحتی کار -۱ و ۱ در نظر گرفتیم. یک ستون ۱ هم اضافه کردم تا محاسبات بایاس مدل راحت تر انجام شود. این تغییرات روی هر دو دسته ی تست و آموزش قبل از کار با دیتاست ها انجام می گیرد. دسته بند هر دو مرحله یک پرسپترون است که نرخ یادگیری آن را برابر ۰,۰۰۱ گرفتیم و شرط توقف الگوریتم تکرار یادگیری و آپدیت وزن ها، یا صفر شدن تعداد داده های اشتباه کلاس بندی شده است و یا حداکثر ۱۰۰ بار تکرار فرایند. که با توجه به درهم رفتگی داده های این مساله شرط اول اتفاق نمی افتد چون داده ها با مدل خطی جداشدنی نیستند.

نمودار دو ویژگی روی شکل زیر داده های قرمز نظیر برچسب -۱ و زرد نظیر برچسب ۱ هستند.



مدل پرسپترون این مرحله یک بردار سه تایی به صورت زیر است.

$$w_0 + w_1 \text{ Feature 1} + w_2 \text{ Feature 2} = 0$$

$$\forall a = [a, b]$$

$$\begin{cases} w_0 + w_1 a + w_2 b \geq 0 & \rightarrow \text{assigned label} = 1 \\ w_0 + w_1 a + w_2 b < 0 & \rightarrow \text{assigned label} = -1 \end{cases}$$

$$\Rightarrow w = [w_0, w_1, w_2]^T$$

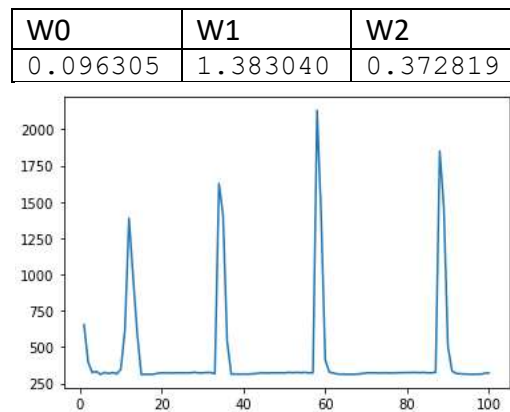
این مقادیر با به روز کردن هرباره ی وزن ها با تمام داده های آموزش (batch mode) از رابطه ی زیر به دست می آید.

$$w_{new} = w_{old} + \eta \sum_{n \in M} x^{(n)} y^{(n)}$$

$$w_{new} = w_{old} + \eta \left[\sum_{n \in M} y^{(n)}, \sum_{n \in M} (bias + y^{(n)}), \sum_{n \in M} bias \cdot y^{(n)} \right]^T$$

این مدل در تابع `train_perceptron` پیاده می شود. هربار با پارامترهای به دست آمده مرحله قبل با تابع `classify` برچسب تمام مجموعه ی آموزش را با رابطه ی XW به دست می آورد و در ستون 'assigned_label' قرار می دهد. سپس تابع `find_missed` با مقایسه ی مقدار برآورد شده با label داده ها ، در ستون 'missclass' تعیین می کند که داده ی مورد نظر درست یا اشتباه دسته بندی شده. و سپس اپدیت وزن ها از رابطه بالا محاسبه می شود.

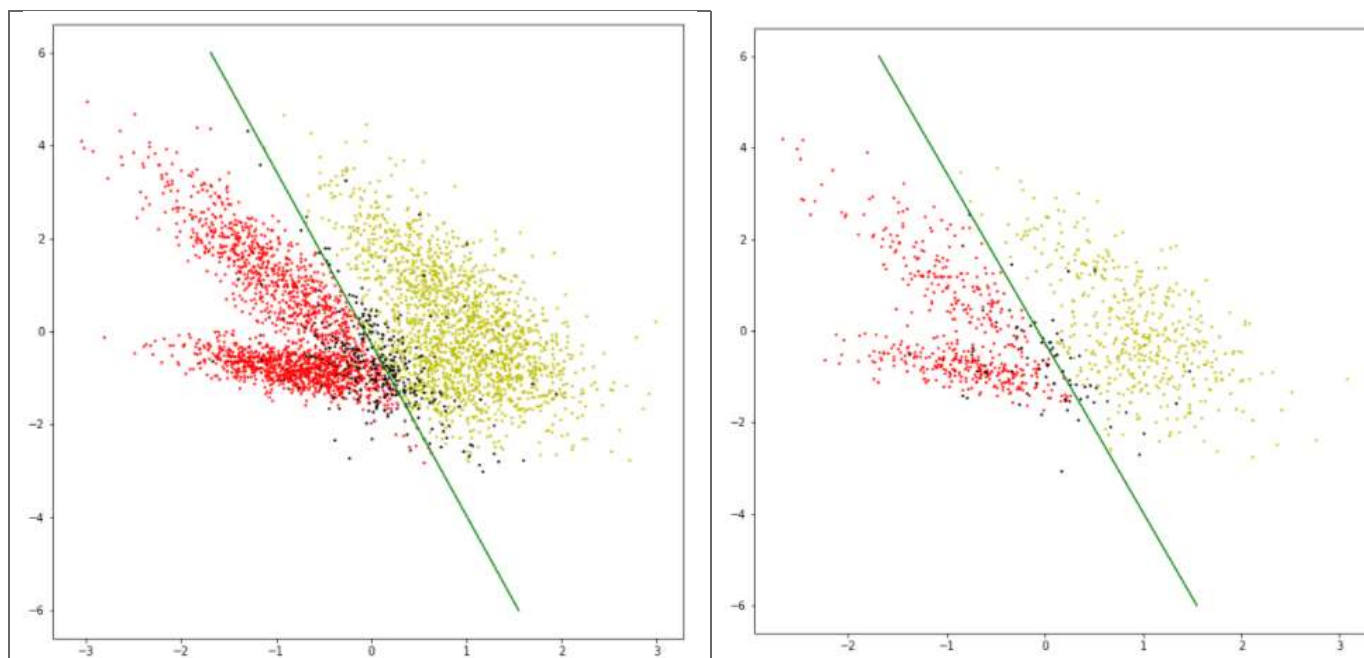
حاصل آموزش داده ها برداری به صورت زیر است. همچنین تعداد خطاهای طبقه بندی روی داده های آموزش در هر مرحله در طول تکرارها در زیر آمده.



الگوریتم اگر در یکی از قله های نمودار متوقف می شد مدل خطای بالا داشت و قابل استفاده نبود. بهتر از شرط توقف ترکیبی از خطا و تعداد تکرار باشد. اینجا مدل من دچار چنین مشکلی نشده پس نتایج را با همین وضعیت می شود تحلیل کرد.

داده های تست و آموزش را با مدل برازش می دهیم و با شمارش تعداد اشتباه کلاس بندی شده ها معیارهای دقت و خطای زیر به دست می آید و همینطور روی نمودارهای هر دو دسته داده های تست و آموزش خط سبز رنگ مدل پرسپترون است و داده های رنگی داده هایی هستند که درست دسته بندی شده اند و نقاط سیاه رنگ نقاطی هستند که مدل نتیجه ی خطا برای نام کلاس آن ها برآورد کرده.

train	test
total = 4000 total_missclassified = 323 misclassified_ratio = 0.08075 correct_ratio = 0.91925	total = 1000 total_missclassified = 73 misclassified_ratio = 0.073 correct_ratio = 0.927



هرچند استفاده از کلاس بندی خطی بهترین مدل برای دسته بندی این داده ها نیست (چون مرز خطی ندارند) اما عملکرد مدل هم روی داده های تست و هم روی داده های آموزشی دقتی بالاتر از ۹۰ درصد دارد. اگر مدل پیچیده تری انتخاب کنم که علاوه بر رابطه ی خطی بین ویژگی ها با برچسب، تاثیر موارد دیگری مثل حاصل ضرب ویژگی ها درهم یا رابطه های غیرخطی از آن ها را با برچسب در نظر بگیرد دقت مدل هم روی داده های تیت هم آموزش بالاتر می رود (تا جایی که دچار خطای بیش برآزش و پیچیدگی شدید مدل نشده باشد). ولی حتی همین مدل خطی هم عملکرد مناسبی دارد.

با استفاده از LDA بردار X که دو ویژگی دارد را به فضایی با یک بعد کمتر می برم. یعنی تمام داده های را به یک عدد تبدیل می کنم و روی این عدد یک مدل پرسپترون برآزش می دهم. این بار مدل پرسپترون من به صورت زیر خواهد بود و دو پارامتر دارد.

$$\begin{aligned}
 w_0 + w_1 \cdot \text{Feature_label} &= 0 \\
 x &= \text{Feature_label} \\
 w_0 + w_1 x &\geq 0 \rightarrow \text{assigned label} = 1 \\
 w_0 + w_1 x &< 0 \rightarrow \text{assigned label} = -1
 \end{aligned}$$

برای محاسبه ی ویژگی LDA نظیر داده ها باید برداری را در راستایی پیدا کنم که اگر داده های دو کلاس روی آن راستا نگاهشته شوند، حداکثر فاصله را از داده های کلاس مخالف و حداقل فاصله را از داده های کلاس مشابه داشته باشند. برای این کار به دنبال یک بردار در فضای دوبعدی دو ویژگی feature1 و feature2 می گردم که تابع زیر را بیشینه کند.

$$\begin{aligned}
 J(w) &= \frac{(\mu'_1 - \mu'_2)^2}{S_1'^2 + S_2'^2}, \quad w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \\
 x &= \begin{bmatrix} \text{feature 1} \\ \text{feature 2} \end{bmatrix} \rightarrow x' = w^T x \in \mathbb{R}
 \end{aligned}$$

که در این تابع :

$$\mu_1' = \frac{1}{N_1} \sum_{y^{(1)}_i=1} x^{(1)}_i = w^T \mu_1, \quad S_1' = \sum_{y^{(1)}_i=1} (x^{(1)}_i - \mu_1')^2$$

$$\mu_2' = \frac{1}{N_2} \sum_{y^{(1)}_i=2} x^{(1)}_i = w^T \mu_2, \quad S_2' = \sum_{y^{(1)}_i=2} (x^{(1)}_i - \mu_2')^2$$

آنچه برای محاسبه ی LDA داده ها نیاز داریم راستای بردار است. پس لازم نیست مقدار بیشینه را حساب کنیم. تنها به دست آوردن راستا کافی است. با توجه به روابط زیر :

$$|\mu_1' - \mu_2'|^2 = |w^T \mu_1 - w^T \mu_2|^2 = w^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T w$$

$$S_1'^2 = \sum_{y^{(1)}_i=1} (w^T x^{(1)}_i - w^T \mu_1)^2 = w^T \left(\sum_{y^{(1)}_i=1} (x^{(1)}_i - \mu_1)(x^{(1)}_i - \mu_1)^T \right) w$$

$$S_2'^2 = \sum_{y^{(1)}_i=2} (w^T x^{(1)}_i - w^T \mu_2)^2 = w^T \left(\sum_{y^{(1)}_i=2} (x^{(1)}_i - \mu_2)(x^{(1)}_i - \mu_2)^T \right) w$$

$$\Rightarrow J(w) = \frac{w^T S_B w}{w^T S_W w} \quad \text{ماتریس بین گروهی کوواریانس}$$

$$S_B + S_W = S \quad \text{ماتریس کوواریانس کل}$$

$$S_B w = \lambda S_W w \quad \text{معادله لاگرانژ}$$

$$\Rightarrow \lambda w = S_W^{-1} S_B w$$

$$\Rightarrow w \propto S_W^{-1} S_B w \quad \text{از دو طرف با } S_W^{-1} \text{ ضرب می کنیم}$$

$$\Rightarrow w \propto S_W^{-1} (\mu_1 - \mu_2)$$

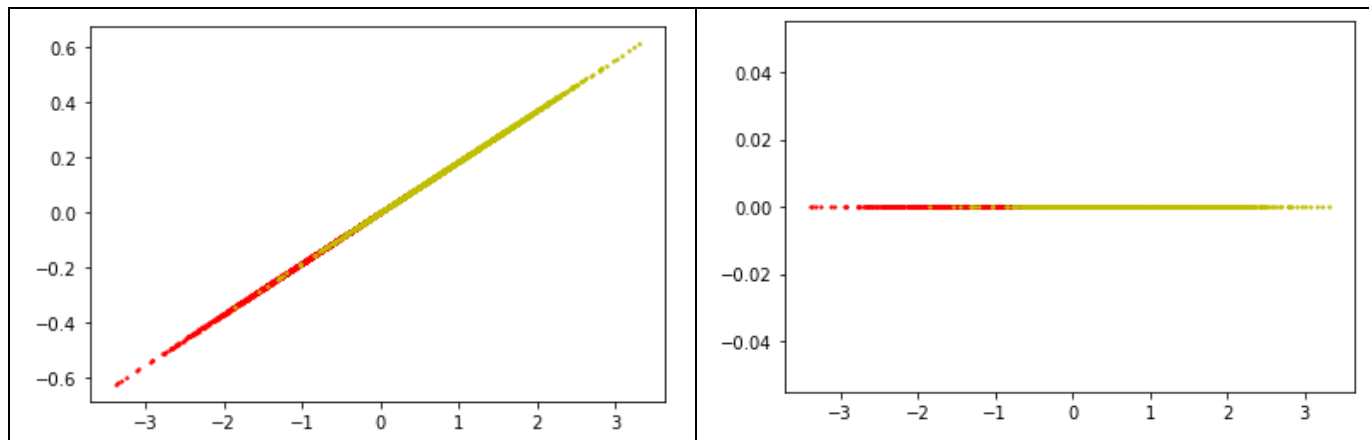
کافیست وارون ماتریس S_W و تفاضل میانگین های دو ویژگی داده های آموزشی را محاسبه کنیم تا راستای w پیدا شود. تابع LDA(d) این کار را پیاده سازی کرده و تمام داده ها را روی این راستا نگاشت می دهد و عدد به دست آمده را در ستونی با نام lda روی داده ها نگه می دارد. این ستون همان ورودی پرسپترون جدید برای ساختن مدل و کلاس بندی خواهد بود. در نهایت یک بردار یکه در راستای مورد نظر برمی گرداند.

قبل از اینکه داده ها را به LDA بدهیم، یک استانداردسازی خواسته شده روی ستون های feature1 و feature2 از رابطه زیر انجام می دهیم. (همین استانداردسازی باید برای داده های تست هم قبل از ارسال به تابع LDA با همین میانگین و واریانس بدست آمده از روی داده های آموزش انجام بگیرد. این میانگین و واریانس داده های آموزش را برآوردی از میانگین و واریانس توزیع اصلی و مجهول داده ها در نظر گرفتیم.)

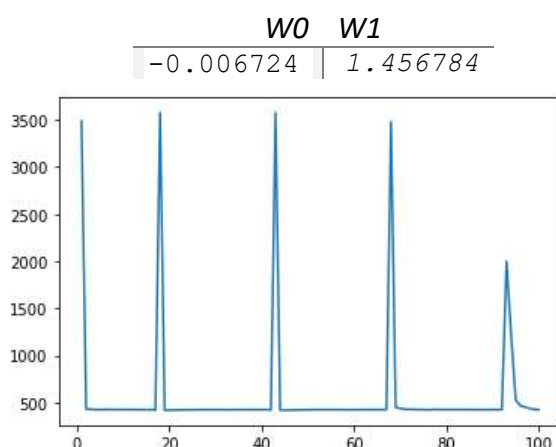
بردار w به دست آمده :

$$\begin{bmatrix} 0.98342694 \\ 0.18130485 \end{bmatrix}$$

مقدار ستون **lda** رو روی خط افقی نشان داده شده میبینیم و اگر این عدد را برای هر داده روی راستای **w** پیدا کنیم، خط به دست آمده شکل زیر همان راستای پیدا شده است که تابع **L** را بیشینه کرد. در زمان کار روی داده های تست و ارزیابی مدل، مقدار **lda** داده های تست با همین بردار محاسبه خواهد شد.



مدل پرسپترون را با همان روشی در قسمت قبل توضیح دادم آموزش می دهیم. همان طور که در روابط بالاتر گفته شد این مدل یک ورودی و دو پارامتر دارد و این ورودی همان **lda** است. مدل به دست آمده :



در مورد قله های نمودار آبی و شرط توقف الگوریتم همچنان آنچه گفته بودم صادق است ولی اینجا چون روی خطای کمی متوقف شده، کار را ادامه می دهیم. بردار مدل پرسپترون به دست آمده بایاس تقریباً صفر دارد. در واقع این مدل چون روی داده های یک بعدی کار می کند که این داده ها از برآزش روی یک بردار یک به دست آمده اند و ابرصفحه ی این فضا یک نقطه ست پس تبدیل می شود به یک مقدار ثابت یا آستانه که یک تابع پله ای طبق آن مقدار بین دو کلاس تمایز قایل می شود. (تابع پله ای اینجا **sign** بوده).

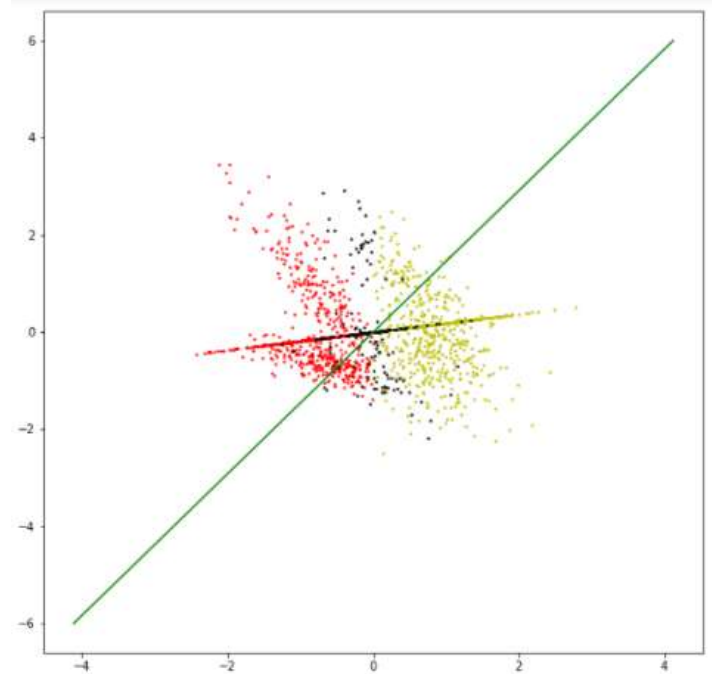
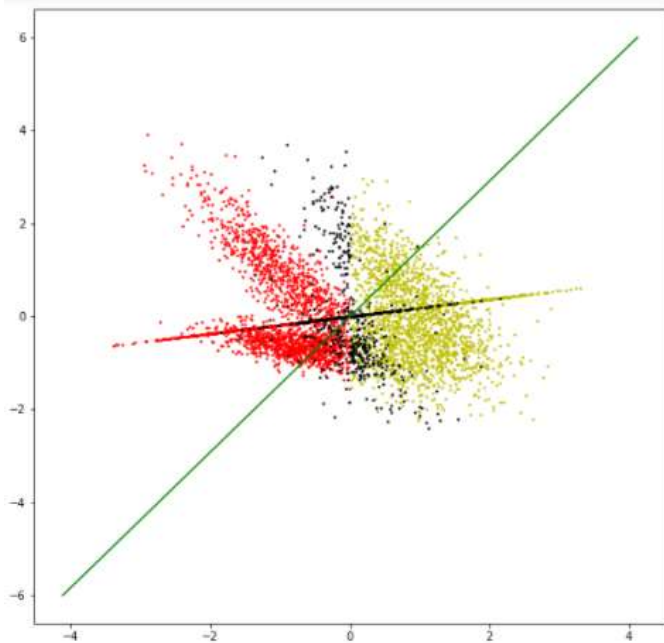
برآزش داده های تست و آموزش نتایج زیر را دارد. روی نمودارهای این بخش، خط سبز مدل پرسپترون است. خط دوم همان نگاشت مقادیر **lda** روی راستای **w** هست. محل تلاقی این دو خط، معیار دیته بندی است. یعنی به ازای داده هایی که سمت چپ آن نقطه مقدار **lda** شان باشد در دسته قرمز ۱- و در غیر این صورت در دسته ی زرد ۱ قرار خواهند گرفت. داده های رنگی داده های درست دسته بندی شده هستند و داده های سیاه اشتباه دسته بندی شده اند. توجه شود که نقاط پراکنده ی نشان داده شده (خارج خط) همان نقاط اصلی **feature1** و **feature2** هستند و فقط برای مشخص شدن درست و غلط کلاس بندی شدنشان نمایش داده شدند تا مقایسه این حالت با حالت بدون **LDA** از روی شکل واضح تر باشد. این مقادیر **feature1** و **feature2** مستقیماً در مدل وارد نشده اند و مدل اینجا یک پرسپترون دو پارامتری است که صرفاً روی داده های خط **lda** عمل می کند. و فضای داده های ورودی تک بعدی هست.

train

test

```
total =
4000
total_missclassified =
428
misclassified_ratio =
0.107
correct_ratio =
0.893
```

```
total =
1000
total_missclassified =
99
misclassified_ratio =
0.099
correct_ratio =
0.901
```



وضعیت کلی دسته بندی روی هر دو سری داده ها به هم بسیار نزدیک است مثل حالت قبل. اگر نمودار خط افقی داده های تست و آموزش را در نظر بگیرم، اط روس این داده های تک نقطه ای که مدل روی آن ساخته شده واضح است که نمیتوان یک نقطه ی آستانه روی آن پیدا کرد. پس داده های یک بعدی LDA این مساله همچنان به حالت خطی جداپذیر نیستند. و روی دو نمودار بالا و نقاط سیاه رنگ، خط LDA را اگر تنها و بدون شیب در نظر بگیرم، همان خط افقی که قبل تر برای داده های آموزش رسم کردم را دارم که به جای برچسب های واقعی، با مقدار های برآورد شده ی مدل رسم شده اند و رنگ سیاه تمام نقاطی ست که کلاس واقعی شان در آستانه ی پیدا شده برای این مدل صدق نمی کند. با وجود این درهم رفتگی داده ها بازهم مدل کارایی چندانی ندارد. ولی اگر نتایج را با نتایج مدل پرسپترون قبلی مقایسه کنم، هرچند اختلاف خیلی زیاد نیست اما افت معناداری در دقت مدل هم در داده های تست و هم آموزش پیدا شده. در بخش قبل گفتم که اصل خطای مدل ناشی از پیچیدگی کم آن است. اتفاقی که افتاد اینکه با تک بعدی کردن داده های مساله، مدل و ورودی هایش بازهم کمی ساده تر شدند و آن جداکنندگی که مدل قبلی قادر بود روی داده های با دو ویژگی بیاید در حالت اسکالر شدن داده ها و با مدل دوم که فقط از یک نقطه تشکیل شده بود محدودتر شد. نتیجه اینکه حدود سه درصد که افت دقت مدل روی داده های تست مدل دوم نسبت به اولی. استفاده از LDA روی این داده ها نتیجه ی مطلوب تر یه دست نمی دهد.

فرض : N داده m کلاس ، n ویژگی و از کلاس $1 \leq j \leq m$ برابر N_j داده وجود دارد .

$$X^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}_{n \times 1} \longrightarrow X'^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_{m-1}^{(i)} \end{bmatrix}_{(m-1) \times 1}, \quad X' = W^T X^{(i)}$$

$$\Rightarrow W = \begin{bmatrix} w_1 & w_2 & \dots & w_{m-1} \end{bmatrix}_{n \times (m-1)}$$

داده‌ها را از فضای n بعدی به فضای $m-1$ بعدی می‌بریم

از قبل برای حالت دو کلاس ماتریس براندگی در دو کلاس :

$$S_w = S_1 + S_2$$

$$S_j = \sum_{y^{(i)}=j} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T, \quad \forall j \in \{1, 2\}$$

حالا تکریف به ازای $\forall j \in \{1, 2, \dots, m\}$ کلاس :
و زک محاسب همان رابطه بالا که در آن زیر برابر است :

$$S_w = \sum_{j=1}^m S_j$$

هر S_j یک ماتریس $n \times n$ است .
 S_w هم یک ماتریس $n \times n$ هست .

$$\mu_j = \frac{1}{N_j} \sum_{y^{(i)}=j} x^{(i)}$$

حالا برای ماتریس براندگی بین کلاسی در حالت دو کلاسی داریم :
یک μ برای کل داده تعریف می‌کنیم .

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

$$\mu = \frac{1}{N} \sum_{i=1}^N x^{(i)}$$

حالا به ازای هر $1 \leq j \leq m$ می‌خواهیم زیرماتریس بین فاصله با μ (مورد کل داده‌ها) داشته باشیم .
از طرفی می‌خواهیم این تأثیر به صورت وزن دار باشد . یعنی هرچه تعداد اعضای کلاس j بیشتر باشد ، داده‌های کلاس j فاصله بیشتری با μ داشته تا جایی بیشتر بین داده‌های X' پیدا شود پس :

$$S'_B = \sum_{j=1}^m N_j (\mu_j - \mu)(\mu_j - \mu)^T$$

که یک ماتریس $n \times n$ هست .

$$J(w) = \frac{w^T S_B w}{w^T S_w w}$$

ما هم مثل حالت در لایه داریم.

که برای این فامده بین داده های در لایه مخفی بهترین و داده های بین لایه های در خوشه ها بهترین براندازی را داشته باشد. $J(w)$ باید بیشترین شود.

$$\max J(w)$$

$$\Rightarrow \frac{\partial J}{\partial w} = 0 \quad \text{در مشتق گیری مثل حالت در لایه ...}$$

$$S_w^{-1} S_B w = \lambda w \quad \hookrightarrow J(w) \in \mathbb{R}$$

در حالت در لایه در لایه هم ...

حالا برای حالت m لایه در w با n هر یک $n-1$ بردار n آبی داریم. برای هر w_j :

$$S_w^{-1} S_B w_j = \lambda_j w_j \quad , \quad j \in \{1, 2, \dots, m-1\}$$

$$J(w_j) \in \mathbb{R}$$

$$w = [w_1, w_2, \dots, w_{m-1}] \quad \begin{matrix} \text{ماتریس} \\ n \times (m-1) \end{matrix}$$

و هر w_j مثل حالت در لایه برابر بردار ویژه $S_w^{-1} S_B$ هست.

اگر هر یک از w_j های نظیر بزرگترین λ_j باشد آمده را بیاکنیم و در w قرار دهیم، مقدار

بهینه $J(w)$ را می توان پیدا کرد و $w = [w_1, \dots, w_j]$ همان ابرصفحه مورد نظر برای

کلاس داده ها است.