

$$\mu = \frac{1}{4} ([1^4] + [0^1] + [1^2] + [2^3]) = \begin{bmatrix} 2.5 \\ 1.5 \end{bmatrix}$$

$$\Sigma = \frac{1}{4} \left(([1^4] - \begin{bmatrix} 2.5 \\ 1.5 \end{bmatrix})([1^4] - \begin{bmatrix} 2.5 \\ 1.5 \end{bmatrix})^T + \dots + ([2^3] - \begin{bmatrix} 2.5 \\ 1.5 \end{bmatrix})([2^3] - \begin{bmatrix} 2.5 \\ 1.5 \end{bmatrix})^T \right)$$

$$= \frac{1}{4} \begin{bmatrix} 5 & 1 \\ 1 & 17 \end{bmatrix} \rightarrow \text{مقادیر ویژه: } (\frac{5}{4} - \lambda)(\frac{17}{4} - \lambda) - 1 = 0 \Rightarrow \lambda^2 - \frac{11}{4}\lambda + \frac{21}{4} = 0$$

$$\rightarrow \lambda_1 = \frac{11 + \sqrt{37}}{4}, \lambda_2 = \frac{11 - \sqrt{37}}{4}$$

$$v_1 = \begin{bmatrix} -4 + \sqrt{37} \\ 1 \end{bmatrix}, v_2 = \begin{bmatrix} -4 - \sqrt{37} \\ 1 \end{bmatrix}$$

حالا دیتاها را در سطرهای یک ماتریس می چینیم، اسر X، خواص داشته

$$X = \begin{bmatrix} 4 & 1 \\ 1 & 0 \\ 2 & 5 \\ 2 & 4 \end{bmatrix} \rightarrow Xv_1 = \begin{bmatrix} 1.22 \\ 0.08 \\ 5.17 \\ 4.25 \end{bmatrix}, Xv_2 = \begin{bmatrix} -47.22 \\ -12.08 \\ -9.17 \\ -32.22 \end{bmatrix}$$

در ~~دیتا~~ ^{راستار} v_1 داده هر کلاس فاصلات زیادی دارند یعنی دایان بالا ولی در ~~دیتا~~ ^{راستار} v_2 فاصلات کمتر دارند \leftarrow در راستار v_2 بیشتر داده ها قابل جوابی اند. همچنین در v_2 نمی توانیم در مقدار یک بعد تفکیک کنیم ولی در v_1 می توانیم یک threshold بنویسیم ۴.۲۵، ۱.۲۳ که دو تادسته جوابی شوند (ب)
 به ازای هر برادر v که به تمامی داده ها اضافه شود، داریم که:

$$\mu' = \sum (x_i + v) = \mu + v$$

$$\rightarrow \Sigma' = \frac{1}{n} \sum_i (x_i - \mu') (x_i - \mu')^T = \frac{1}{n} \sum_i (x_i + v - \mu - v) (x_i + v - \mu - v)^T = \frac{1}{n} \sum_i (x_i - \mu) (x_i - \mu)^T = \Sigma$$

$\rightarrow \Sigma' = \Sigma \rightarrow$ pc! don't change

$$f_{\text{ensemble}}(x) = \frac{1}{M} \sum_{i=1}^M f_i(x)$$

برای سادگی؛ f_e نشان می‌دهیم

الف)

$$\text{Var}(f_e) = \text{Var}\left(\frac{1}{M} \sum_{i=1}^M f_i(x)\right) = \frac{1}{M^2} \text{Var}\left(\sum_{i=1}^M f_i(x)\right)$$

سوال
طبق فرض f_i ها مستقل بودند پس داریم که

$$\frac{1}{M^2} \sum_{i=1}^M \text{Var}(f_i(x))$$

طبق فرض سوال $\text{Var}(f_i)$ ها برابر اند
پس داریم که

$$\frac{1}{M^2} \times M \text{Var}(f_1(x)) = \frac{\text{Var}(f_1(x))}{M}$$

$$\Rightarrow \boxed{\text{Var}(f_e) = \frac{\text{Var}(f_1(x))}{M}}$$

حالا اگر M را زیاد کنیم، واریانس f_e کم خواهد شد.

حال bias را بررسی می‌کنیم.

$$\text{Bias} = E[\hat{f}(x)] - f(x)$$

تابع واقعی \rightarrow تابع تخمینی

f_e که ثابت است پس $E[f_e(x)]$ را حساب می‌کنیم.

چون f_i ها بایس یکسانی دارند طبق فرض سوال.

$$E[f_e(x)] = E\left[\frac{1}{M} \sum_{i=1}^M f_i(x)\right] = \frac{1}{M} \sum_{i=1}^M E(f_i(x)) = E(f_1(x))$$

$$\rightarrow \text{Bias} = E(f_1(x)) - f(x) \rightarrow$$

ϕ بایس با تخمین M تخمینی نمی‌کند و ثابت می‌ماند.

ب) ابتدا واریانس را در این حالت حساب می‌کنیم، داریم که

$$\text{Var}\left(\frac{1}{M} \sum_{i=1}^M f_i(x)\right) = \frac{1}{M^2} \text{Var}\left(\sum_{i=1}^M f_i(x)\right) = \frac{1}{M^2} \left(\sum_{i=1}^M \text{Var}(f_i(x)) + 2 \sum_{1 \leq i < j \leq M} \text{Cov}(f_i(x), f_j(x)) \right)$$

طبقی فرض سوال داریم $\text{cov}(f_i(x), f_j(x)) = \rho$

$$\Rightarrow \text{var}(f_c) = \frac{1}{m^2} \left(m \text{var}(f_1(x)) + \rho \left(\frac{m^2 - m}{1} \right) \right) = \frac{\text{var}(f_1(x))}{m} + \frac{(m^2 - m)}{m^2} \rho$$

$$\Rightarrow \left[\text{var}(f_c) = \frac{\text{var}(f_1(x))}{m} + \frac{m-1}{m} \rho \right] = \rho + \frac{1}{m} \left(\text{var}(f_1(x)) - \rho \right)$$

\downarrow
 $(1 - \frac{1}{m}) \rho$

با افزایش m به سمت ∞ ، مقدار واریانس برابر با ρ می شود

حالا اگر $\rho > 0$ در صورت ثابت بودن m ، با افزایش ρ ، واریانس زیاد می شود

اگر $\text{var}(f_1(x)) - \rho > 0$ باشد ، با افزایش m واریانس کمتر می شود

اگر $\text{var}(f_1(x)) - \rho = 0$ باشد ، با افزایش m واریانس تغییر نمی کند

اگر $\text{var}(f_1(x)) - \rho < 0$ باشد ، با افزایش m واریانس زیاد می شود

در مورد بایاس ، چون E خطای انت و داخلین جمع مرتبتر ا دل داریم ، تغییر ایی دنی شود نسبت به حالت مستقل \leftarrow بایاس باز هر تغییر نمی کند

با یادگیرنده های ضعیف ada boost نیاز به مستقی پذیر بودن ندارند همانطور که در الگوریتم موجود در اسلایدها داریم ، آمیت زن های یادگیرنده ضعیف بدون مستقی گیر انجام می شود

boosting هزینه بیشتری دارد ، زیرا boosting به صورت سری کار می کند و معاری انجام نمی شوند

همچنین مدل های پیچیده تر نیز دارند چون از تمام دیتا برای یادگیر استفاده می کنند برخلاف bagging که وی بعضی از دیتا کار می کند

الف) اگر $k=1$ باشد، در نقطه هر تردید ترین به خودش است، پس خطای training صفر می شود.

ب) اگر k زیاد باشد، عملاً دایره های بهینه عموم بر حسب هاچی هستند بر اساس این دلیل می زنیم، به عبارتی bias خیلی زیادی داریم. اگر هر k خیلی کم باشد، ممکن است دچار کمی به اشتباه بر حسب داشته باشند تردید آن نقطه باشد به آن نقطه هر اشتباه بر حسب زده می شود، در صورتی که اگر همسایه ها بیشتر می دسیم ممکن بود درست بر حسب بزنیم.

پ) با استفاده از کد یا ستون برابر ۱۲ — $k=1$ ، بهترین k ، ۷۵ بود و با خطای ۰/۲۸۵۷.

ت) در شکل منحنی به آورده شده است.


```

import numpy as np
from sklearn.model_selection import LeaveOneOut, cross_val_score
from sklearn.neighbors import KNeighborsClassifier

pos = np.array([[2, 7], [3, 8], [5, 1], [6, 2], [7, 3], [8, 4], [9, 5]])
neg = np.array([[1, 5], [2, 6], [3, 7], [4, 8], [5, 9], [7, 2], [8, 3]])
X = np.vstack((pos, neg))
y = np.array([1] * 7 + [0] * 7)
LeaveOne = LeaveOneOut()
errors = {}

for k in range(1, 14):
    knn = KNeighborsClassifier(n_neighbors=k)
    scores = cross_val_score(knn, X, y, cv=LeaveOne, scoring="accuracy")
    errors[k] = 1 - np.mean(scores)

```

errors

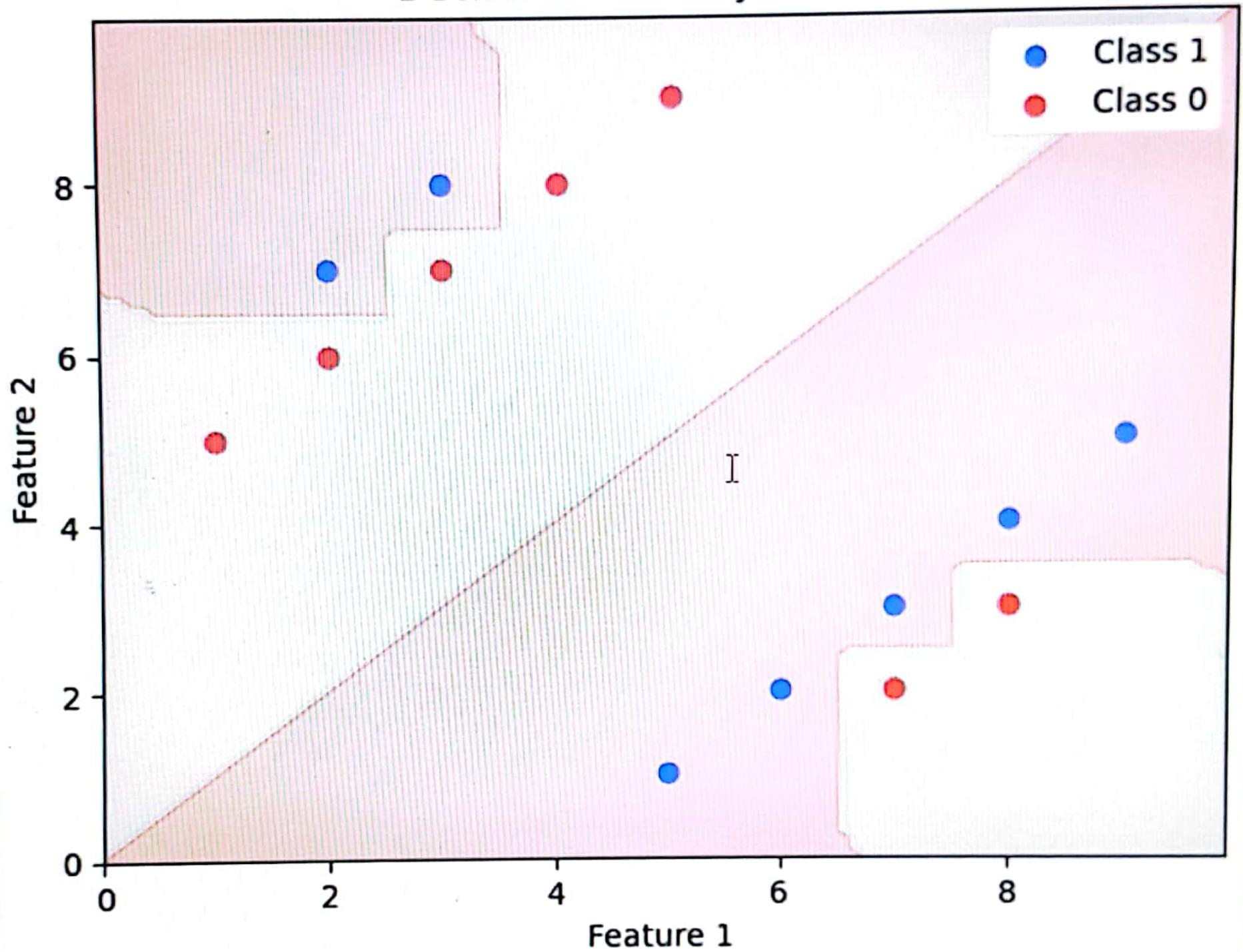
✓ 0.2s

```

{1: 0.7142857142857143,
 2: 0.7142857142857143,
 3: 0.4285714285714286,
 4: 0.6428571428571428,
 5: 0.2857142857142857,
 6: 0.3571428571428571,
 7: 0.2857142857142857,
 8: 0.6428571428571428,
 9: 1.0,
10: 0.7142857142857143,
11: 1.0,
12: 0.7142857142857143,
13: 1.0}

```

Decision Boundary for 1-NN



الف) برای پیدا کردن مینیمم، مشتق می‌گیریم و برابر صفر می‌نویسیم.

~~دستی که~~ $\frac{\partial L}{\partial \mu_t} = \frac{\partial}{\partial \mu_t} \sum_{j=1}^k \sum_{x_i \in S_i} \|x_i - \mu_t\|^2$

دستی که $j=t$ است مشتق ناهمراست ←

$$\frac{\partial L}{\partial \mu_t} = \frac{\partial}{\partial \mu_t} \sum_{x_i \in S_t} \|x_i - \mu_t\|^2 \Rightarrow -2 \sum_{x_i \in S_t} (x_i - \mu_t) = 0$$

$$\rightarrow n\mu_t - \sum_{x_i \in S_t} x_i = 0 \rightarrow \boxed{\mu_t^* = \frac{\sum_{x_i \in S_t} x_i}{n}}$$

همان مقدار میانگین نقاط هر خوشه می‌شود μ_t یعنی برابر آن خوشه.

ب) • الگوریتم همواره همگرا می‌شود زیرا در هر مرحله داریم خطا را کمتر می‌کنیم (طبق اصلادهها) و تعداد حالات انقباض مرکز ثقل محدود است ← الگوریتم همگرا می‌شود علت کاهش خطا هم این است که هر مرحله داریم در جهت عکس گام‌ها μ حرکت می‌کنیم.

الگوریتم به مقدار کمی اولیه خوشه‌ها حساس است و برابر بهبود انقباض اولیه مرکز ثقل، الگوریتم‌هایی ارائه شده است. (ممکن است چند تا مرکز ثقل داشته باشیم و GD در مینیمم محلی بنشیند)

ب) اگر z داخل خوشه‌ای قبلی نباشد، ممکن است هم بین خوشه‌ها قرار بگیرد و همگرای الگوریتم دچار اختلال کند در این صورت ~~این روند~~ هر سری خطا کم شود، برقرار نمی‌شود.

۵- بار پیدا کردن \argmin عبارت $\|x - \alpha u\|^2$ از این عبارت نسبت به α مشتق می‌گیریم و برابر صفر قرار می‌دهیم.

$$\frac{\partial}{\partial \alpha} (x - \alpha u)^T (x - \alpha u) = -2u^T (x - \alpha u) = 0 \rightarrow u^T x - \alpha u^T u = 0 \rightarrow \boxed{\alpha = u^T x}$$

$$\Rightarrow \boxed{f_u(x) = \frac{x^T u}{u^T u} u}$$

بردار \hat{v} ای بوده $\alpha = u^T x$ و $\hat{v} = \alpha u$ بود.

$$\rightarrow \argmin_{\substack{u \\ u^T u = 1}} \frac{1}{m} \sum \|x^{(i)} - f_u(x^{(i)})\|^2 = \frac{1}{m} \sum \|x^{(i)} - u^T x^{(i)} u\|^2$$

$$= \frac{1}{m} \sum (x^{(i)T} x^{(i)} + (u^T x^{(i)})^2 - 2(u^T x^{(i)})^2)$$

$$= \frac{1}{m} \sum (x^{(i)T} x^{(i)} - (u^T x^{(i)})^2) \rightarrow \argmin_{\substack{u \\ u^T u = 1}} \text{عبارت} \equiv$$

$$* \argmax_{\substack{u \\ u^T u = 1}} \frac{1}{m} \sum (u^T x^{(i)})^2 \rightarrow \argmax_u \|Xu\|^2$$

ماتریک X که به صورت $\begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix}$ بردار دیده می‌شود متناظر با بزرگترین مقدار ویژه است، در این صورت می‌کند.

طبق اسلایدها داریم که:

$$\text{Var}(XV) = \frac{1}{n} \sum_{i=1}^n (x_i^T V)^2$$

$$\rightarrow \text{Var}(XV) = \frac{1}{n} \|XV\|^2 = \frac{1}{n} V^T X^T X V = \frac{1}{n} V^T \Sigma V$$

عبارت $\frac{1}{n} V^T \Sigma V$ هم معادل همین رابطه است، در PCA هر V بردار ویژه متناظر با بزرگترین مقدار ویژه است.

ماتریک MSE بین نقاط و تقویر آنها را کمینه می‌کند، همان مولفه اول PCA است.