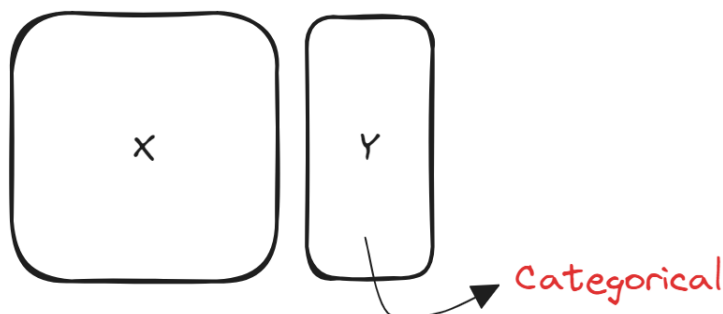


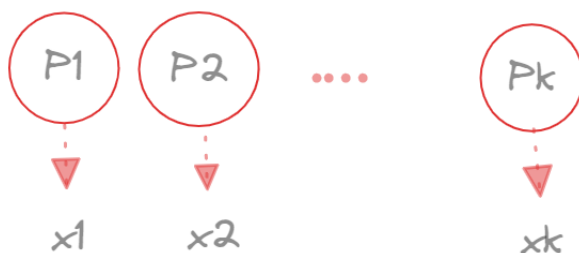
## دسته بندی (Classification)

در اینجا، داده (که به شکل ماتریس  $X$  نشان داده می شود) و وکتور نهایی ( $Y$ ) را داریم. این بار پاسخ ها به صورت Categorical هستند.



شکل ۱: ماتریس ورودی و خروجی

یکی از راه هایی که برای حل این نوع مسائل داریم، محاسبه کردن  $P(X)$  اول تا  $k$  ام و در نظر گرفتن بزرگ ترین آن هاست. به این روش Maximum Likelihood می گوئیم.



شکل ۲: احتمال هر یک از دسته ها

$$h : X \rightarrow \{1, 2, \dots, k\}$$

$$P_1 \rightarrow x \xrightarrow{h(x)} P_2$$

$$E = \mathbb{1}\{Y \neq h(X)\}$$

این فرمول نشان دهنده مقدار ارور ما می باشد که به ازای هر اشتباه ۱ واحد به ارور اضافه می شود. (در اینجا  $Y$  و  $X$  هر دو متغیر تصادفی هستند.)

$$P(Y, X) = P(X|Y)P(Y)$$

$$P_1(X) = P(X|Y = 1)$$

$$P_2(X) = P(X|Y = 2)$$

هدف این است که تابع تخمینی را پیدا کنیم که امید ریاضی خطا را کمینه کند.

$$\min_h \mathbb{E}[E]$$

که در آن

$$h : X \rightarrow \{1, 2, \dots, k\}$$

تابع فرضیه است.

## MAP - Maximum A Posteriori

در آمار بیزی، MAP به معنای حداکثر احتمال پسین است. این تکنیک راهی برای تخمین پارامترهای یک توزیع احتمال پسین است که با استفاده از داده‌های مشاهده شده و توزیع احتمال پیشین، به دست می‌آید.

۱. توزیع احتمال پسین (Posterior Distribution)

$$P(\theta|D) \propto P(D|\theta) \cdot P(\theta)$$

۲. Map Estimation

$$\theta_{MAP} = \arg \max_{\theta} P(\theta|D)$$

نمادهای فوق، معانی زیر را می‌رسانند.

- $P(\theta|D)$ : توزیع احتمال پسین پارامترها ( $\theta$ ) با توجه به داده‌های مشاهده شده ( $D$ )
- $P(D|\theta)$ : احتمال داده‌ها به شرط پارامترها.
- $P(\theta)$ : توزیع احتمال پیشین پارامترها.

از MAP می‌توان برای به دست آوردن یک تخمین نقطه‌ای از یک کمیت مشاهده نشده بر اساس داده‌های تجربی استفاده کرد.

$$\operatorname{argmax} P(Y|X) = \max_{i \in \{1, 2, \dots, k\}} P(Y = i|X)$$

$$\operatorname{argmax} P(Y|X) = \frac{P(X|Y = i)P(Y = i)}{P(X)}$$

و چون  $\operatorname{argmax}$  را روی  $i$  می گیریم پس باید صورت کسر را بیشینه کنیم.

$$\operatorname{argmax} P(X|Y = i)P(Y = i)$$

به این روش دسته بندی، روش بیز گفته می شود و به هر دسته بند که به این روش عمل کند، Bayes Classifier می گوئیم.

اگر  $P_i$  ها را نداشته باشیم، می توانیم:

(۱) از نمونه ها استفاده کرده و  $P_i$  ها را تخمین بزنیم.  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

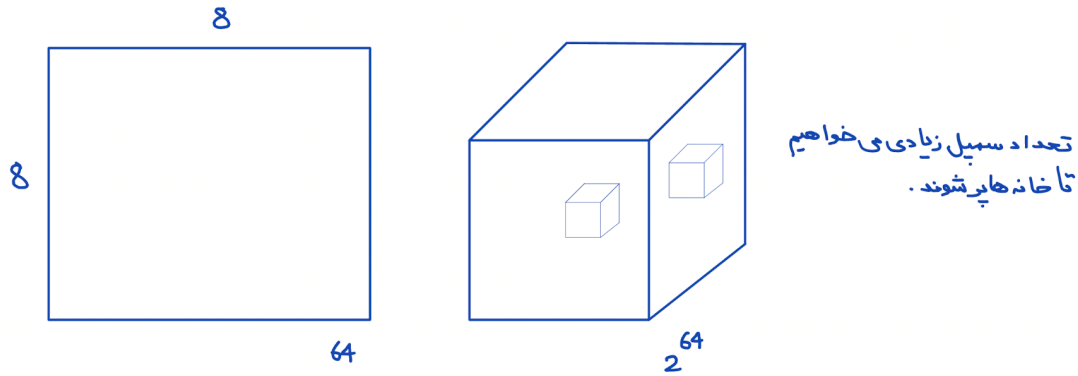
به این روش **Density Estimation** می گوئیم. (۲)

$$\begin{pmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_k(x) \end{pmatrix} \left( \xrightarrow{\text{Softmax}} \right) \begin{pmatrix} \frac{e^{h_1(x)}}{z} \\ \frac{e^{h_2(x)}}{z} \\ \vdots \\ \frac{e^{h_k(x)}}{z} \end{pmatrix}$$

$$z = \sum e_i^h(x)$$

در Density Estimation اولین کاری که باید انجام دهیم، بدست آوردن هیستوگرام از نمونه های موجود است. توزیع نمونه ها در هر قسمت به صورت Binomial می باشد.

در این جا با چالش هایی رو به رو هستیم که یکی از آن ها نفرین ابعاد یا Curse of Dimensionality می باشد.



نفرین ابعاد یا Curse of Dimensionality یک مفهوم در زمینه‌ی آمار و یادگیری ماشین است و به چالش‌ها و مشکلاتی اشاره دارد که در فضاهای با ابعاد بالا (تعداد زیادی از ویژگی‌ها یا متغیرها) به وجود می‌آید. برخی از چالش‌ها و مسائل مرتبط با Curse of Dimensionality عبارتند از:

- تنوع ویژگی‌ها: با افزایش ابعاد فضا، تنوع (تفاوت) بین نقاط داده‌ها نسبت به همدیگر افزایش می‌یابد. این امر ممکن است منجر به افزایش فاصله‌های اقلیدسی بین نقاط شود.
- نیاز به داده‌های بیشتر: با افزایش ابعاد، نیاز به تعداد داده‌های آموزشی بیشتر برای کسب یک مدل قابل اعتماد نیز افزایش می‌یابد. این امر می‌تواند منجر به مشکل کمبود داده در بعضی مسائل شود.
- داده‌های Sparse: در فضاهای با ابعاد بالا، داده‌ها اغلب به صورت sparse می‌شوند؛ به عبارت دیگر، اکثر نقاط در فضا خالی هستند. این موضوع باعث افزایش پیچیدگی در مدل‌سازی و تحلیل داده می‌شود.
- پیدا کردن الگوها و روابط معنادار: افزایش ابعاد ممکن است باعث کاهش قدرت تفسیر و تحلیل الگوها و روابط معنادار در داده‌ها شود.

برای مقابله با Curse of Dimensionality، رویکردهایی نظیر انتخاب ویژگی، تجمیع داده، و الگوریتم‌های خاص مخصوص فضاهای با ابعاد بالا مورد استفاده قرار می‌گیرد. همچنین، استفاده از روش‌های کاهش ابعاد مانند تحلیل مؤلفه‌های اصلی (PCA) و تحلیل تفسیری مؤلفه‌ها (ICA) نیز به عنوان راه‌حل‌های معمول به شمار می‌آید.

## Naive Bayes

$$P(x_1, x_2, \dots, x_p)$$

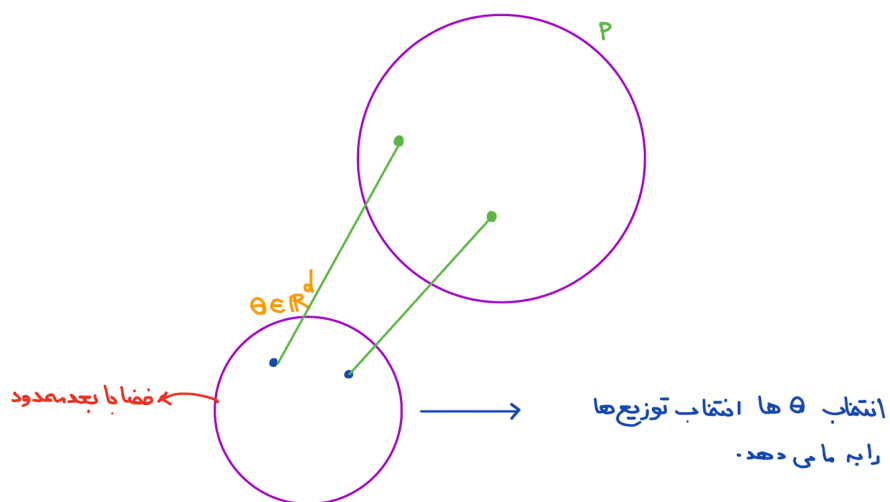
$\Downarrow$

$$P(x_1)P(x_2) \dots P(x_p)$$

$\Downarrow$

$$P(x_1)P(x_2|x_1)P(x_3|x_2, x_1) \dots P(x_p|x_{p-1}, \dots, x_1)$$

$$P \begin{cases} \text{Parametric} \\ \text{Non-parametric} \end{cases}$$



اگر دو  $\theta$  یک توزیع را به ما بدهد، می گوئیم این توزیع نسبت به  $\theta$  ها Identifiable نیست.

## Parametric

$$X \sim P_{\theta} \xrightarrow{\text{Density}} \lambda e^{-\lambda x}$$

$$X \sim P_{\theta} \xrightarrow{\text{Dist}} 1 - e^{-\lambda x}$$

## Non Parametric

$$\text{Density } f(x) \quad , \quad \int f''(x)^2 dx \leq \text{Threshold}$$

## LDA: Linear Discriminant Analysis

فرض کنیم دو توزیع توأم گاوسی با احتمال پیشین  $\pi$  داریم. یعنی:

$$\pi_1, \mathcal{N}(\mu_1, \Sigma) \longrightarrow \theta_1 = (\hat{\pi}_1, \hat{\mu}_1, \hat{\Sigma})$$

$$\pi_2, \mathcal{N}(\mu_2, \Sigma) \longrightarrow \theta_2 = (\hat{\pi}_2, \hat{\mu}_2, \hat{\Sigma})$$

قرار است  $\theta_1$  و  $\theta_2$  را تخمین بزنیم. می‌توانیم از Maximum Likelihood استفاده کنیم اما با مشکلاتی مواجه می‌شویم.. مشکل اصلی ما این است که  $\Sigma$ ، از دو توزیع مختلف آمده است و ۲ تخمین متفاوت برای آن بدست می‌آید و همچنین ابعاد آن به صورت quadratic زیاد می‌شود. بنابراین به دنبال راه حل دیگری می‌رویم.

### راه حل

فرض کنیم داده‌ها به صورت  $x \in X \subseteq \mathbb{R}^d$  از دو توزیع گاوسی آمده‌اند.

$$f_x = (x|y = 1) = \mathcal{N}(\mu_1, \Sigma_1)$$

$$f_x = (x|y = 0) = \mathcal{N}(\mu_2, \Sigma_2)$$

$$p(y = 1) = \pi_1$$

و هدف ما این است که جهت نامساوی را در عبارت زیر بدست آوریم:

$$p(y = 0|x) \stackrel{?}{\leq} p(y = 1|x)$$

سعی می‌کنیم نامساوی بالا را حل کنیم.

$$g(x) = \frac{p(y = 1|x)}{p(y = 0|x)} \leq 1$$

$$\longrightarrow \frac{\frac{p(x|y=1)p(y=1)}{p(x)}}{\frac{p(x|y=0)p(y=0)}{p(x)}} \leq 1$$

$$\longrightarrow \frac{\pi_1 \frac{1}{\sqrt{(2\pi)^d |\Sigma_1|}} \exp\left(-\frac{1}{2}((x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1))\right)}{\pi_2 \frac{1}{\sqrt{(2\pi)^d |\Sigma_2|}} \exp\left(-\frac{1}{2}((x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2))\right)} \leq 1$$

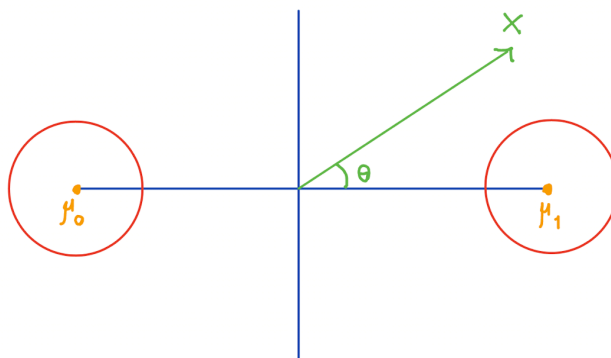
از دو طرف لگاریتم گرفته و به صورت زیر رابطه را بازنویسی می‌کنیم.

$$g(x) = \log\left(\frac{\pi_1}{\pi_2}\right) + \frac{1}{2} \log\left(\frac{|\Sigma_2|}{|\Sigma_1|}\right) + \frac{1}{2}((x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)) - \frac{1}{2}((x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)) \leq 0$$

تابع  $g(x)$  یک تابع درجه دو است.  $\Sigma_*$  و  $\Sigma_1$  را یکی در نظر می‌گیریم تا تابع به شکل خطی بدست آید.

$$g(x) = \log\left(\frac{\pi_1}{\pi_*}\right) + \left(x - \frac{\mu_1 - \mu_*}{2}\right)^T \Sigma^{-1}(\mu_1 - \mu_*)$$

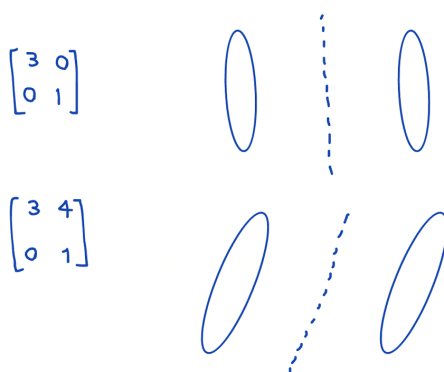
برای مثال اگر  $\pi_1 = \frac{1}{4}$  و  $\Sigma = I_d$



در شکل بالا مکان خط جدا کننده به مقدار  $\pi$  ربط دارد مثلاً اگر  $\pi_1 = \frac{1}{4}$  باشد خط به گروه ۱ نزدیکتر می‌شود.

همچنین شکل توزیع‌ها، جهت خط و زاویه‌دار بودن آن نیز به میزان  $\Sigma$  بستگی دارد و با تغییر آن دوران خواهیم

داشت. به عنوان مثال اگر  $\Sigma = \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}$  باشد، توزیع‌ها به شکل بیضی درمی‌آیند. و اگر  $\Sigma = \begin{pmatrix} 3 & 4 \\ 0 & 1 \end{pmatrix}$  باشد شکل دسته‌بند ما به صورت زیر خواهد بود.



حالا فرض کنیم توزیع دو کلاس گاوسی با  $\Sigma$  برابر هستند. پارامترهای زیر را می‌خواهیم با استفاده از MLE تخمین بزنیم.

$$\hat{\mu}_* = \frac{1}{n_*} \sum_{y_i=0} x_i$$

$$\hat{\mu}_1 = \frac{1}{n_1} \sum_{y_i=1} x_i$$

$$\hat{\Sigma} = \frac{1}{n-2} \sum_{i=1}^n \sum_j (x_j^i - \hat{\mu}_i)(x_j^i - \hat{\mu}_i)^T$$

$$\hat{\pi} = \frac{n_1}{n}$$

اینجا می‌توان  $\Sigma$  ها را یکی در نظر نگرفت و به صورت

$$\Sigma_1 = \Sigma + y_1 I, \Sigma_2 = \Sigma + y_2 I$$

نوشت. در این صورت دیگر تابع خطی نمی‌شود ولی با این وجود، به سختی درجه دو هم نمی‌باشد.