



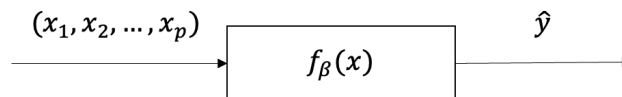
## یادگیری ماشین

نیم سال دوم ۱۴۰۱-۱۴۰۲

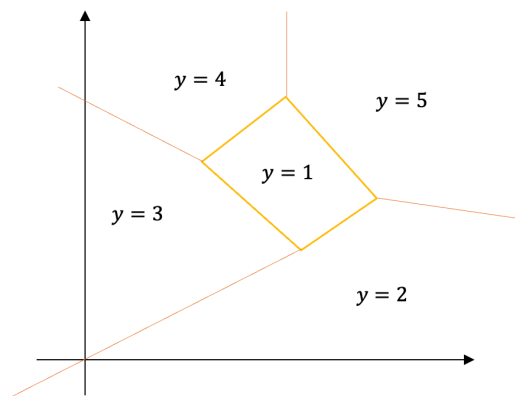
مدرس: دکتر سید ابوالفضل مطهری

## درسنامه ششم

در مسأله دسته‌بندی می‌خواهیم از روی ورودی یکی از  $k$  دسته مختلف را تخمین بزنیم.



که در آن  $\hat{y} = \{1, 2, \dots, k\}$  می‌باشد. آن چیز که مشخص است، تابع  $f_\beta(x)$  فضای ورودی را به  $k$  قسمت افراز می‌نماید و هر افراز به یکی از دسته‌ها تخصیص می‌یابد.



همانطور که در مسأله رگرسیون بحث گردید در اینجا هم نیاز به داشتن یک تابع ریسک و یک مدل هستیم تا بتوانیم تعمیم خارج داده را بدست آوریم.

## دسته‌بند بیز

فرض نمایید که توزیع مشترک  $x$  و  $y$  موجود است یعنی  $P(x, y)$  را در اختیار داریم. می‌دانیم MAP (Maximum A Posteriori) کمترین خطا را برای تصمیم‌گیری از روی ورودی  $x$  در اختیار می‌گذارد.

$$\hat{y} = \arg \max_k \mathbb{P}(y = k | x)$$

برای یک مسأله دو کلاسه کافیت که نسبت احتمالات فوق را با یک مقایسه کنیم:

$$\frac{\mathbb{P}(y = 1 | x)}{\mathbb{P}(y = 2 | x)} \geq 1$$

و اگر از طرفین لگاریتم بگیریم خواهیم داشت:

$$\log\left(\frac{\mathbb{P}(y = 1|x)}{\mathbb{P}(y = 2|x)}\right) \geq 0$$

دسته‌بند فوق را دسته‌بند بیز می‌نامیم و این دسته‌بند تابع زیر را کمینه می‌نماید:

$$\mathbb{P}(y \neq \hat{y}) = \mathbb{E}[\mathbb{I}_{y \neq \hat{y}}]$$

از روی رابطه فوق یکی از توابع ریسکی که می‌توانیم از روی داده بدست آوریم به صورت زیر خواهد بود.

$$\mathcal{L}(y, \hat{y}) = \mathbb{I}_{y \neq \hat{y}}$$

## (LDA) Linear Discriminant Analysis

یک مسأله  $k$  کلاسه را در نظر بگیرید که توزیع مشترک  $x$  و  $y$  به صورت زیر تعریف شده است:

$$\begin{aligned} \mathbb{P}(y = 1) &= \pi_1; & \mathbb{P}(x|y = 1) &= f_1(x) = \mathcal{N}(\mu_1, \Sigma_1) \\ &\vdots & &\vdots \\ \mathbb{P}(y = k) &= \pi_k; & \mathbb{P}(x|y = k) &= f_k(x) = \mathcal{N}(\mu_k, \Sigma_k) \end{aligned}$$

با توجه به دسته‌بند بیز کافیت رابطه زیر را محاسبه نماییم:

$$\begin{aligned} \hat{y} &= \arg \max_k \mathbb{P}[y = k|x] \\ &= \arg \max_k \mathbb{P}[x|y = k] \mathbb{P}[y = k] \\ &= \arg \max_k \frac{\pi_k e^{-\frac{(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)}{2}}} {|\Sigma_k|^{1/2}} \\ &= \arg \max_k \log(\pi_k) - \frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \\ &= \arg \max_k \delta_k(x) \end{aligned}$$

به  $\delta_k(x)$  تابع discriminant می‌گوییم.

در حالت اول فرض نماییم که تنها دو کلاس داریم و  $\Sigma = \Sigma_1 = \Sigma_2$  در این حالت:

$$\begin{aligned} \delta_1(x) &= \log(\pi_1) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \\ \delta_2(x) &= \log(\pi_2) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (x - \mu_2)^T \Sigma^{-1} (x - \mu_2) \end{aligned}$$

مرز بین دو ناحیه را می‌توانیم از روی تساوی قرار دادن دو عبارت فوق بدست آوریم:

$$\delta_1(x) = \delta_2(x) \\ \Rightarrow \log\left(\frac{\pi_1}{\pi_2}\right) + (\mu_1 - \mu_2)^T \Sigma^{-1} x - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 = 0$$

رابطه فوق نشان می‌دهد که ناحیه تصمیم‌گیری یک ناحیه خطی است. در حقیقت که ابرصفحه وجود دارد که بین کلاسهای ۱ و ۲ تفکیک می‌گذارد. این ابرصفحه دارای بردار نرمال

$$w^T = (\mu_1 - \mu_2)^T \Sigma^{-1}$$

می‌باشد و به صورت زیر قابل بیان است:

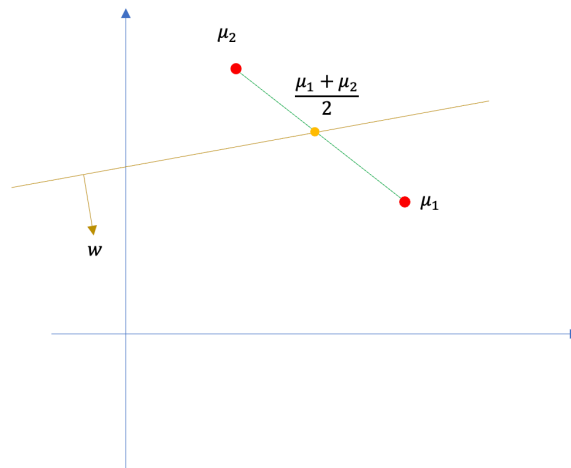
$$w^T x + b = 0$$

که در آن

$$b = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \log\left(\frac{\pi_1}{\pi_2}\right)$$

به طور خاص‌تر اگر  $\pi_1 = \pi_2 = \frac{1}{2}$  باشد آنگاه نقطه  $\frac{\mu_1 + \mu_2}{2}$  بر روی ابرصفحه قرار می‌گیرد. یعنی

$$(\mu_1 - \mu_2)^T \Sigma^{-1} \frac{\mu_1 + \mu_2}{2} = \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2$$



اگر  $\Sigma = \mathbb{I}$  باشد آنگاه  $w = \mu_1 - \mu_2$  خواهد بود و ابرصفحه عمود می‌گردد.

به راحتی می‌توان شرایط فوق را تعمیم داد و برای  $k$  دسته دید که نواحی تصمیم‌گیری با یکسری ابرصفحه از یکدیگر جدا می‌گردند.

در یک مسأله یادگیری داده‌ها در دسترس می‌باشند و توزیع را در اختیار نداریم. با توجه به آنکه در مدل فوق توزیعهای هر دسته گوسی مفروض است از روی داده‌ها می‌توانیم پارامترهای توزیع را تخمین بزنیم. در این شرایط

$$\hat{\pi}_k = \frac{N_k}{N}$$

که در آن  $N_k$  تعداد مشاهدات از کلاس  $k$  می باشد.

$$\hat{\mu}_k = \frac{\sum_{y_i=k} x_i}{N_k} = \frac{\sum x_i I_{y_i=k}}{\sum I_{y_i=k}}$$

$$\hat{\Sigma} = \frac{\sum_{k=1}^K \sum_{y_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T}{N - K}$$

## مقابله با داده‌های نامتقارن

در پاره‌ای از موارد تعداد نمونه‌های یک کلاس کم می باشد و بنابراین دسته‌بند نمی تواند به خوبی برای آنها پاسخگو باشد. چند راه حل می توان برای این مسأله در نظر گرفت:

۱. **نمونه برداری (Resampling):** در این روش یا به دسته کم از طریق نمونه برداری تصادفی از مشاهدات، نمونه های تکراری اضافه می کنیم و یا اینکه به طور تصادفی از دسته های بزرگتر نمونه دور می ریزیم. در این روش ممکن است اطلاعات مفید از بین رفته و یا اینکه overfitting رخ دهد.

۲. **تولید نمونه های تصادفی:** از روی نمونه های با تعداد کم، یک مدل یاد گرفته و داده های جدید مصنوعی تولید می کنیم. از این روش با نام SMOTE نیز یاد می شود که مخفف Synthetic Minority Over-Sampling Technique است.

۳. **روش های Ensemble:** از جمله این روش ها می توان به Bagging و Boosting اشاره کرد که در این درس نامه به آنها نمی پردازیم.

۴. **یادگیری بر پایه متعادل کردن خطا:** خطای کلاسه بندی را می توان به دو خطای مختلف تفکیک نمود:

$$\mathbb{P}[\hat{y} \neq y] = \mathbb{P}[\hat{y} \neq 1 | y = 1] \mathbb{P}[y = 1] + \mathbb{P}[\hat{y} \neq 2 | y = 2] \mathbb{P}[y = 2]$$

اگر توزیع مشترک  $y$  و  $\hat{y}$  را در یک جدول نمایش دهیم، داریم:

		$y$	
		$\pi_1$	$\pi_2$
$\hat{y}$		1	2
	$\pi'_1$	$P_{11}$	$P_{12}$
	$\pi'_2$	$P_{21}$	$P_{22}$

در این شرایط می توان نوشت:

$$\mathbb{P}[\hat{y} \neq 1 | y = 1] = \frac{P_{21}}{P_{11} + P_{21}} = \alpha_1$$

$$\mathbb{P}[\hat{y} \neq 2 | y = 2] = \frac{P_{12}}{P_{12} + P_{22}} = \alpha_2$$

اگر  $y = 1$  را با علامت  $-$  یعنی بیمار نبودن طرف و  $y = 2$  را با علامت  $+$  یعنی بیمار بودن طرف نشان دهیم، آنگاه از عبارات زیر استفاده می‌کنیم:

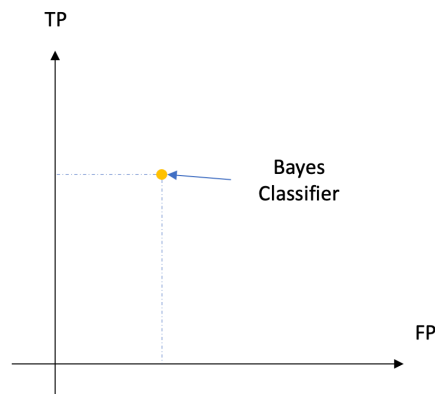
$\alpha_2$  : False Negative (FN)

$\alpha_1$  : False Positive (FP)

$1 - \alpha_1$  : True Negative (TN)

$1 - \alpha_2$  : True Positive (TP)

معمولا تعداد مثبت‌ها در جامعه بیشتر از منفی‌ها می‌باشد و بنابراین تشخیص آن‌ها مهم است. در این شرایط به TP حساسیت (sensitivity) هم می‌گوییم. همچنین به TN که تشخیص گروه بزرگتر است specificity می‌توانیم FP و FN بیانگر حساسیت و specificity می‌باشد. می‌توانیم TP را در یک صفحه نمایش دهیم:



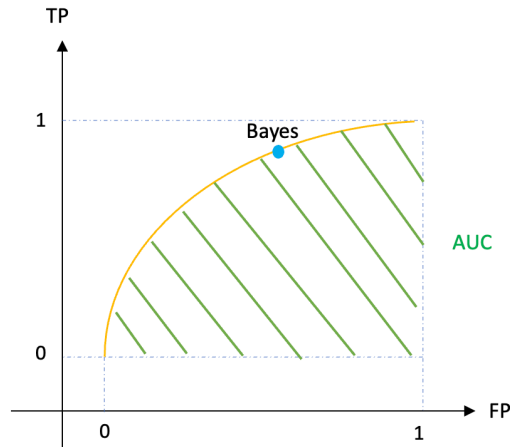
مشاهده کردیم که دسته‌بند Bayes میزان

$$\mathbb{P}[y \neq \hat{y}] = \alpha_1 \pi_1 + \alpha_2 \pi_2 = (FP)\pi_1 + (1 - TP)\pi_2$$

را کمینه می‌کند. در عمل دوست داریم  $FP \rightarrow 0$  و  $TP \rightarrow 1$  میل کند. قضیه Neyman-Pearson بیان می‌دارد که برای مصالحه بین  $(\alpha_2, \alpha_1)$  و یا  $(FP, TP)$  کفایت مقدار نسبت احتمالات را با یک آستانه مقایسه نماییم. یعنی:

$$\log \frac{\mathbb{P}(y = 1|x)}{\mathbb{P}(y = 2|x)} \geq Th$$

و به ازای آستانه‌های مختلف TP و FP-های مختلف می‌گیریم. در نتیجه شکل نمودار به صورت زیر خواهد بود:



بنابراین با تغییر در آستانه می‌توانیم منحنی فوق را که ROC (Receiver Operating Characteristics) نامیده می‌شود به دست آوریم. مساحت زیر ROC را AUC (Area Under the Curve) می‌نامیم. در حالت ایده‌آل  $AUC = 1$  می‌باشد. یعنی در  $TP = 1$  و  $FP = 0$  قابل حصول است. همه مطالب فوق را برای شرایطی مطرح نمودیم که توزیع در اختیار است. علی‌القاعده برای هر سیستم دسته‌بندی می‌توانیم از روی داده‌های آزمون ماتریس confusion را بدست آوریم که در حقیقت همان ماتریس توزیع مشترک  $(y, \hat{y})$  است. با این تفاوت که به جای فرکانس، تعداد نشان داده می‌شود. از روی این ماتریس می‌توانیم همه پارامترها را تخمین بزنیم. همچنین با تغییر آستانه در LDA می‌توانیم ROC و AUC را نیز بدست آوریم و از کیفیت کار مطلع شویم.

## Quadratic Discriminant Analysis (QDA)

روابط LDA را با فرض  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_k$  بدست آوردیم. حال اگر این فرض را کنار بگذاریم، آنگاه تابع discriminant به شکل زیر می‌شود.

$$\sigma_k(x) = -\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1}(x - \mu_k) - \frac{1}{2} \log |\Sigma_k| + \log \pi_k$$

در این شرایط نواحی تصمیم‌گیری با توابع درجه دوم از یکدیگر جدا می‌گردند.

