

# FEATURE SELECTION

در فصل‌های گذشته، فرض کردیم ویژگی‌ها از قبل انتخاب شده‌اند و تنها به طراحی طبقه‌بند پرداخته شد. در این فصل به بررسی روش‌های انتخاب این ویژگی‌ها می‌پردازیم.

- هدف‌ها:
  - ۱- انتخاب تعداد نمونه ویژگی‌ها (l)
  - ۲- انتخاب بهترین l ویژگی

l بزرگ ۳ عیب دارد:

- ۱- پیچیدگی محاسباتی بیشتر
- ۲- کارایی تعیین‌پذیری پایین

تعداد نمونه‌های آموزشی

۳- تخمین خطای ضعیف (برای تخمین خطای خوب باید  $\frac{N}{l}$  زیاد باشد)

معمولاً  $\frac{N}{l} = 10$  یا  $\frac{N}{l} = 20$

تعداد ویژگی‌ها (l)

برای N داده شده: l باید به حدی بزرگ باشد که بتواند الف تفاوت‌های کلیدی را از بین برد. آموزش دهد.

ب) شباهت‌های داخلی هر مدل را آموزش بدهد.

l باید به حدی کوچک باشد که تفاوت‌های نمونه‌های داخلی مدل را آموزش ندهد.

در عمل  $\frac{N}{l} < 10$  یک گزینه قابل قبول برای برخی موارد است.

# FEATURE SELECTION

- ❖ The goals:
  - Select the “optimum” number  $l$  of features
  - Select the “best”  $l$  features
- ❖ Large  $l$  has a three-fold disadvantage:
  - High computational demands
  - Low generalization performance
  - Poor error estimates

➤ Given  $N$

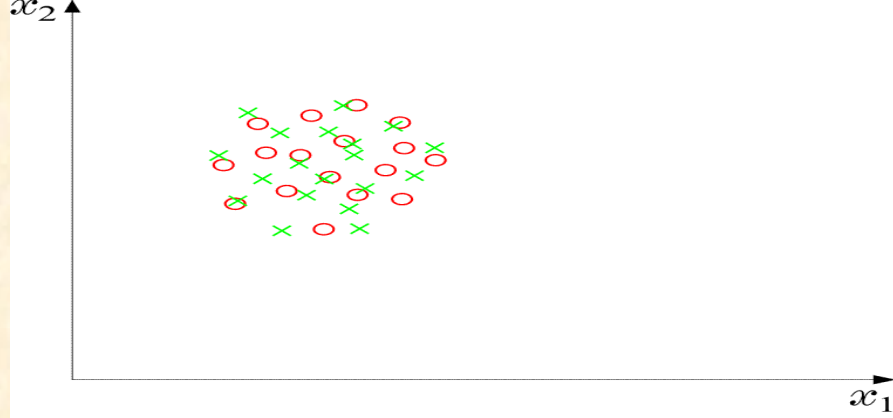
- $l$  must be **large enough** to learn
  - what makes classes **different**
  - what makes patterns in the same class **similar**
- $l$  must be **small enough** not to learn what makes patterns of the same class **different**
- In practice,  $l < N/3$  has been reported to be a sensible choice for a number of cases

کدام  $l$  تا دقت بیشتری دارد؟

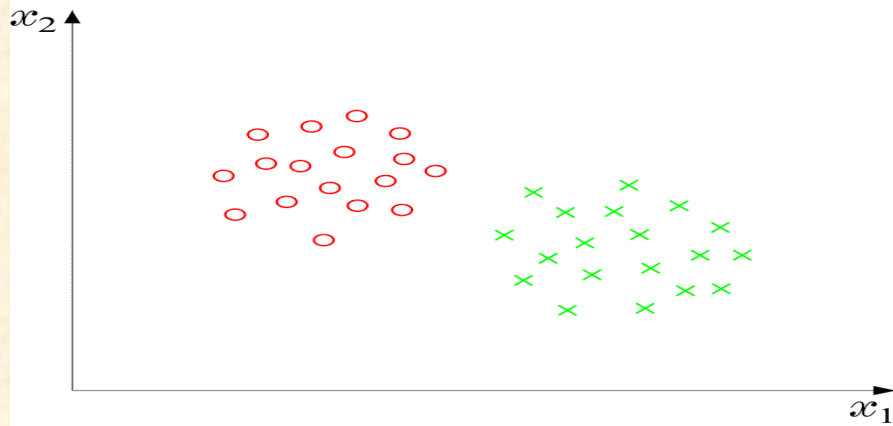
➤ Once  $l$  has been decided, choose the  $l$  most informative features

- Best: **Large** between class distance,  
**Small** within class variance

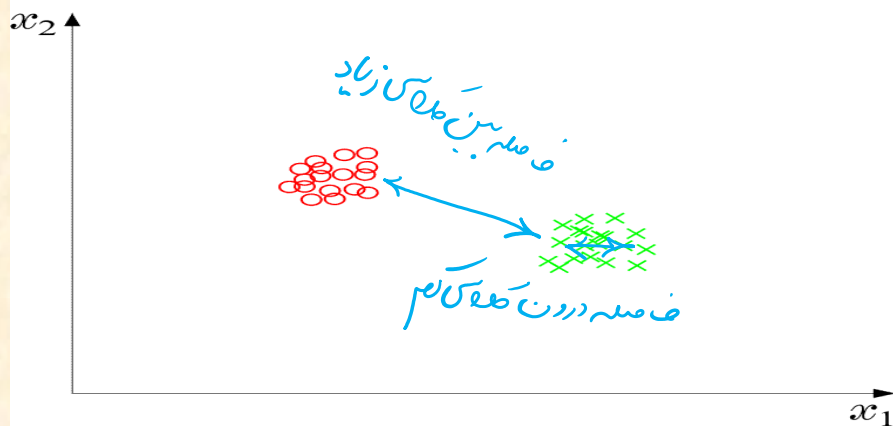
بهترین  $l$  دقت دارد:   
بیشترین فاصله بین کلاس‌ها   
کمترین فاصله درون کلاس‌ها (واریانس داخلی کم)



**Bad choice**



**Not bad choice**



**Good choice**





## 5.2.2 نرمالیزه کردن داده‌ها Data Normalization

وقتی مقادیر ویژگی‌ها در هم‌ردیف‌های مختلفی قرار داده می‌شوند نرمالیزه کردن داده‌ها است.

different dynamic range

زیرا در غیر این صورت ویژگی‌های با مقادیر بزرگتر اثر بیشتری روی نتیجه فرآیند خواهد داشت که ممکن است آن‌ها را حذف کند و در نتیجه از اهمیت آن ویژگی‌ها نباشد.

روش کورات نرمالیزه کردن با استفاده از کمترین و بیشترین دواریابی:

$$\left\{ \begin{array}{l} \bar{x}_k = \frac{1}{N} \sum_{z=1}^N x_{zk} \quad ; \quad k=1, \dots, d \\ \sigma_k^2 = \frac{1}{N-1} \sum_{z=1}^N (x_{zk} - \bar{x}_{zk})^2 \end{array} \right.$$

دوره نرمالیزه شده با بیشترین صفر و بیشترین

روش حذفی

$[0, 1]$  ,  $[-1, 1]$

$$\hat{x}_{zk} = \frac{x_{zk} - \bar{x}_{zk}}{\sigma_k}$$

در کنار روش‌های خطی، روش‌های غیرخطی نیز وجود دارند که برای حالت‌هایی که تفویج داده‌ها اطراف مبدأ منبسط به صورت بی‌نهایت نیست کاربرد دارند. در روش‌های غیرخطی تبدیلی‌هایی با توابع غیرخطی (نظیر گزینشی،  $\log$  و  $\sin$ ) برای تبدیل داده‌ها در بازه‌های مورد نظر به بازه‌های مورد نیاز صورت می‌گیرد.

برای نشان روش مقیاس‌بندی Soft max :  
Scaling

$$y = \frac{x_{ijk} - \bar{x}_{jk}}{r_{jk}}$$

$$\hat{x}_{ijk} = \frac{1}{1 + e^{-y}}$$

- داده‌ها را در محدوده  $[0, 1]$  قرار می‌دهند.

- برای  $y$  های خیلی کوچک رابطه خطی است.

### 5.2.3 : داده‌های گزینشی، رفته Missing Data :

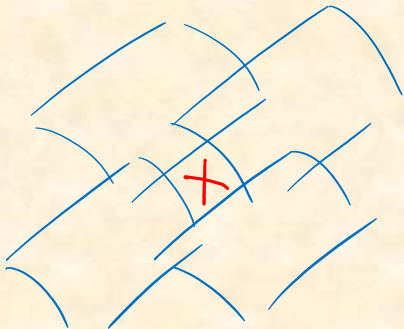
در عمل، برخی ویژگی‌ها از بزرگ داده‌ها گزینشی هستند.

در علوم انسانی: نظرسنجی‌های گزینشی survey : پاسخ‌های ناقص ← داده‌های گزینشی رفته

در پزشکی از دور: نواحی مشخص تنها با بررسی از جغرافیای تحت پوشش قرار می‌گیرد.

در شبکه‌های حسگر: اطلاعات از حسگرها تفویج شده جمع‌آوری و در آنجا ارسال می‌شوند

← داده‌های ناقص



# Imputation

روشهای تکمیل داده‌های ناقص:



۱- قرار دادن صفر به جای داده‌های گمشده است

۲- قرار دادن میانگین بدون شرط که از رول و فیلتر موجود در داده‌ها استفاده می‌شود.

۳- قرار دادن میانگین شرطی، در صورت موجود بودن نمونه‌ها از  $pdf$  های داده‌های گمشده است

فلسفہ اصلی انتخاب دہرے : حذف دہرے ہاں تکی نہ اطہعات ضلیفی رلہرہ .

## ❖ The basic philosophy

- ہر کی تو آ دہرے ہاں قوی ہائی نہ .

- Discard individual features with **poor** information content
- The remaining information rich features are examined **jointly** as vectors

5.4: انتخاب دہرے ہر اسل روش تست فرضیہ آماری :

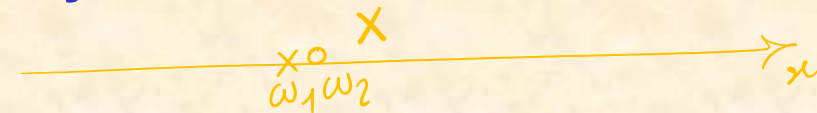
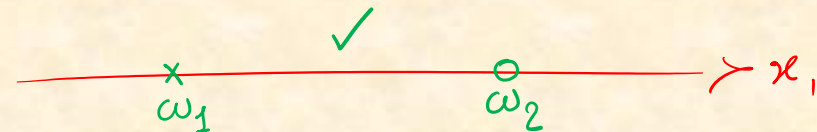
## ❖ Feature Selection based on statistical Hypothesis Testing

- The Goal: For each individual feature, find whether the values, which the feature takes for **the different classes**, **differ significantly**.

That is, answer

- $H_1 : \theta_1 \neq \theta_0$ : The values differ significantly
- $H_0 : \theta_1 = \theta_0$ : The values do not differ significantly

If they do not differ significantly reject feature from subsequent stages.



## ❖ Hypothesis Testing Basics



## 5.4.1 اصول فرضیه آماری :

فرض کنید  $x$  متغیر تصادفی با توزیع احتمال مشخص است که پارامتر  $\theta$  آن نامعلوم است.  
 برای مثال اگر توزیع  $\theta$  دس باشد این پارامتر می تواند میانگین یا واریانس باشد.

$$\begin{cases} H_1: \theta \neq \theta_0 \\ H_0: \theta = \theta_0 \end{cases}$$

روال تقسیم: فرض کنید  $n = 1, 2, \dots, \infty$  نمونه های آموزشی متغیر تصادفی  $x$  باشند  $f(x_1, \dots, x_n)$

انتخاب می شود. بنابراین به منظور  $f(x_1, x_2, \dots, x_n)$  ،  $q$  صدق انتخاب می شود که توزیع

آماره تست test statistics

احتمال  $q$  به راحتی بر حسب پارامتر نامعلوم  $\theta$  فرمول بندی شود:  $P_q(\theta)$

$D$ : بازه ای که در آن با احتمال زیاد فرض  $H_0$  برقرار باشد (بازه قبولی) acceptance interval

$\bar{D}$ : ممکن  $D$  است. فرض  $H_1$  برقرار است. (بازه بحرانی) critical



احتمال تقسیم بندی:

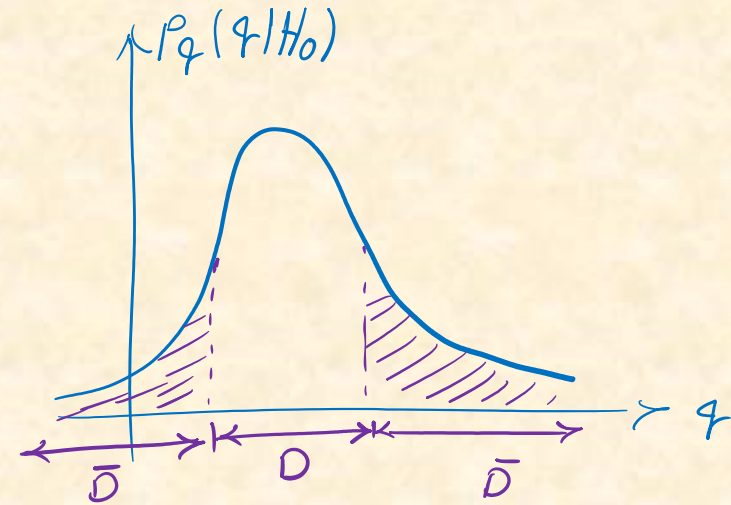
فرض کنید  $H_0$  فرض درست باشد. بنابراین احتمال رخ دادن خطا:

$$p(q \in \bar{D} | H_0) = \alpha$$

$\bar{D}$  می  $P_q(q | H_0)$  انحراف

significance level:  $\alpha$

سطح اهمیت





## حالت با دارایی ملغوا :

فرض کنید  $x$  متغیر تصادفی و  $N = 1, 2, \dots$  نمونه های حاصل از این فرآیند باشند. فرض کنید نمونه ها به صورت  
 متغیر مستقل از هم  $E[x] = \mu$  و  $E[(x - \mu)^2] = \sigma^2$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

تخمین مستقل برای  $\mu$  بر حسب نمونه ها :

از آنجایی که با تغییر نمونه ها تخمین متفاوتی حاصل می شود.  $\bar{x}$  نیز متغیر تصادفی است. به بافتن احتمال  $f_{\bar{x}}(\bar{x})$

$$E[\bar{x}] = \frac{1}{N} \sum_{i=1}^N E[x_i] = \frac{1}{N} N \mu = \mu$$

مشغول می شود :

$\bar{x}$  یک تخمین ناهایب Unbiased از میانگین  $\mu$  است.

$$E[(\bar{x} - \mu)^2] = E\left[\left(\frac{1}{N} \sum_{i=1}^N x_i - \mu\right)^2\right] = \frac{1}{N^2} \sum_{i=1}^N E[(x_i - \mu)^2] +$$

$$\frac{2}{N^2} \sum_{i \neq j} E[(x_i - \mu)(x_j - \mu)]$$

درجه  $N$  بزرگتر باشد، دایس تر تخمین حول  $\mu$  می شود  
 یعنی  $\sigma_{\bar{x}}^2 = \frac{1}{N} \sigma^2$  از فرض استقلال  
 $E[(x_i - \mu)(x_j - \mu)] = 0$

همه چیز در  $N$  بزرگتر می شود

فرض کنید  $\hat{\mu}$  داده شده، باید تصمیم زیر را بگیریم:

$$H_1: E[X] \neq \hat{\mu}$$

$$H_0: E[X] = \hat{\mu}$$

برای این منظور باره تست زیر را تعریف می‌کنیم:

$$q = \frac{\bar{x} - \hat{\mu}}{\sigma / \sqrt{N}}$$

با توجه به قضیه همگونی باین (تفویح احتمال  $\bar{x}$  تحت فرض  $H_0$ ) یعنی  $\hat{\mu}$  (داده شده) تقریباً گادسی است:

$$p_{\bar{x}}(\bar{x}) = \frac{\sqrt{N}}{\sqrt{2\pi}\sigma} e^{-\frac{N(\bar{x} - \hat{\mu})^2}{2\sigma^2}}$$

Appendix A

$$\sim N\left(\hat{\mu}, \frac{\sigma^2}{N}\right)$$

بنابر این تفویح  $q$  تحت فرض  $H_0$  به صورت  $N(0, 1)$

برای سطح اهمیت  $\alpha$ : بازه اعتماد  $D = [-\alpha, \alpha]$  صدق این شرط در این بازه احتمال  $q$

Table 5.1

$1-\alpha$	0.8	0.85	0.9	...	$\alpha$
	1.282	1.440	1.645	...	

( $\alpha$  احتمال حسود در بازه  $0$ )

کامها:

۱- با نمونه آزمودنی داده شده از متغیر  $X$ ،  $q$  را به دست می آوریم.

۲- سطح اهمیت  $\alpha$  را انتخاب می کنیم.

۳- از روی جدول مقادیر  $N(0,1)$ ، بازه قبول  $D = [a, b]$  را انتخاب می کنیم (۱- $\alpha$ ) به دست می آوریم.

۴- اگر  $q \in D$  باشد،  $H_0$  را قبول می کنیم، در غیر این صورت  $H_1$  را قبول می کنیم.

مثال: متغیر تصادفی  $X$  با  $\sigma^2 = 0.23$ ،  $N = 16$ ،  $\bar{x} = 1.35$ ،  $p = 0.05$ ، فرض کنید  $\mu = 1.4$

درست است یا خیر؟

$$\alpha = 0.05 \rightarrow x_p = 1.967 \quad \text{Table 5.1}$$

$$\rightarrow \text{prob} \left\{ -1.97 < \frac{\bar{x} - \hat{\mu}}{\sigma / \sqrt{N}} < 1.97 \right\} = 0.95$$

$$\frac{1.35 - 1.4}{\sqrt{0.23/4}} = \frac{-0.05}{0.0575} = -0.87 \in D$$

$H_0$  قابل قبول  
 $E[X] = \hat{\mu}$

➤ The steps:

- $N$  measurements  $x_i, i = 1, 2, \dots, N$  are known

- Define a function of them

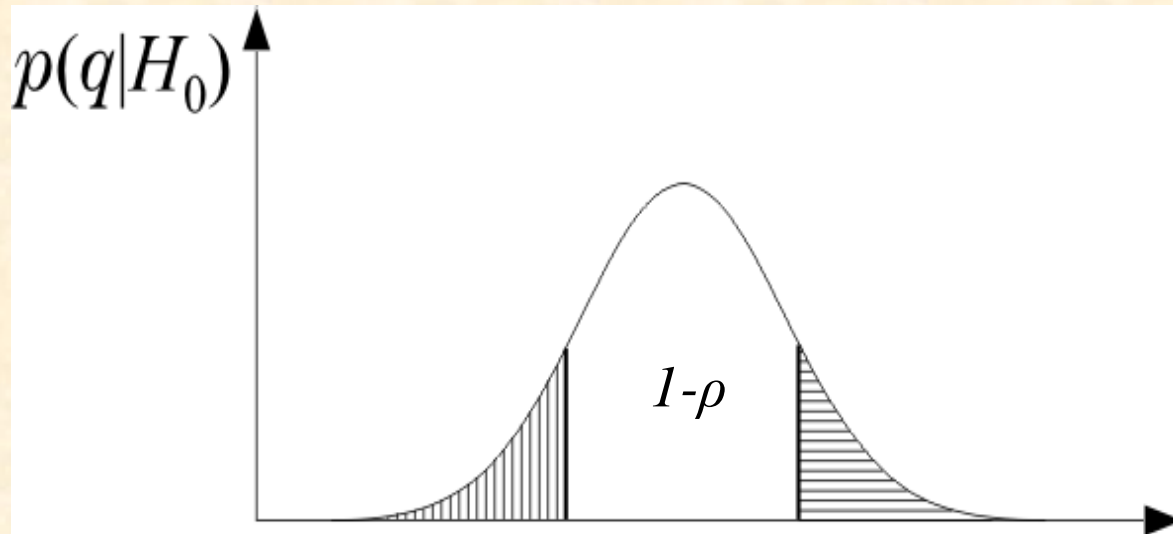
$$q = f(x_1, x_2, \dots, x_N): \quad \text{test statistic}$$

so that  $p_q(q; \theta)$  is easily parameterized in terms of  $\theta$ .

- Let  $D$  be an interval, where  $q$  has a high probability to lie under  $H_0$ , i.e.,  $p_q(q|\theta_0)$
- Let  $\bar{D}$  be the complement of  $D$   
 $D \longrightarrow$  Acceptance Interval  
 $\bar{D} \longrightarrow$  Critical Interval
- If  $q$ , resulting from  $x_1, x_2, \dots, x_N$ , lies in  $D$  we accept  $H_0$ , otherwise we reject it.

➤ Probability of an error

$$p_q(q \in \bar{D} | H_0) = \rho$$



- $\rho$  is preselected and it is known as the **significance level**.

❖ Application: The known variance case:

- Let  $x$  be a random variable and the experimental samples,  $x_i = 1, 2, \dots, N$ , are assumed mutually **independent**. Also let

$$E[x] = \mu$$

$$E[(x - \mu)^2] = \sigma^2$$

- Compute the sample mean

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

- This is also a random variable with mean value

$$E[\bar{x}] = \frac{1}{N} \sum_{i=1}^N E[x_i] = \mu$$

That is, it is an **Unbiased Estimator**

➤ The variance  $\sigma_{\bar{x}}^2$

$$\begin{aligned} E[(\bar{x} - \mu)^2] &= E\left[\left(\frac{1}{N} \sum_{i=1}^N x_i - \mu\right)^2\right] \\ &= \frac{1}{N^2} \sum_{i=1}^N E[(x_i - \mu)^2] + \frac{1}{N^2} \sum_i \sum_j E[(x_i - \mu)(x_j - \mu)] \end{aligned}$$

Due to independence

$$\sigma_{\bar{x}}^2 = \frac{1}{N} \sigma_x^2$$

That is, it is **Asymptotically Efficient**

➤ **Hypothesis test**

$$H_1 : E[x] \neq \hat{\mu}$$

$$H_0 : E[x] = \hat{\mu}$$

➤ **Test Statistic: Define the variable**

$$q = \frac{\bar{x} - \hat{\mu}}{\sigma / \sqrt{N}}$$

- Central limit theorem under  $H_0$

$$p_{\bar{x}}(\bar{x}) = \frac{\sqrt{N}}{\sqrt{2\pi}\sigma} \exp\left(-\frac{N(\bar{x} - \hat{\mu})^2}{2\sigma^2}\right)$$

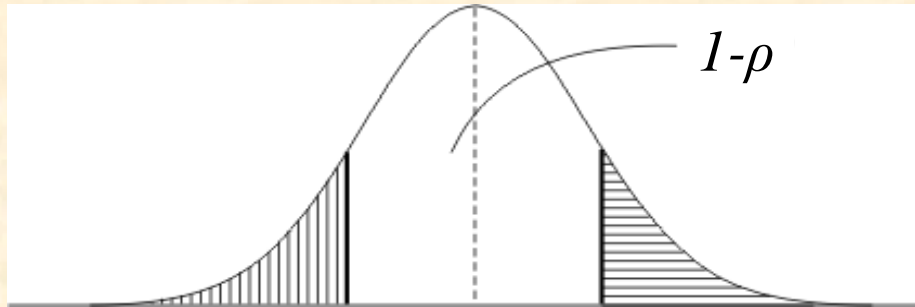
- Thus, under  $H_0$

$$p_q(q) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{q^2}{2}\right) \quad q \approx N(0,1)$$



➤ The decision steps

- Compute  $q$  from  $x_i, i=1,2,\dots,N$
- Choose significance level  $\rho$
- Compute from  $N(0,1)$  tables  $D=[-x_\rho, x_\rho]$



- if  $q \in D$  accept  $H_0$   
if  $q \in \bar{D}$  reject  $H_0$

➤ **An example:** A random variable  $x$  has variance  $\sigma^2=(0.23)^2$ .  $N=16$  measurements are obtained giving  $\bar{x}=1.35$ . The significance level is  $\rho=0.05$ .

Test the hypothesis

$$H_0 : \mu = \hat{\mu} = 1.4$$

$$H_1 : \mu \neq \hat{\mu}$$

➤ Since  $\sigma^2$  is known,  $q = \frac{\bar{x} - \hat{\mu}}{\sigma/4}$  is  $N(0,1)$ .

From tables, we obtain the values with acceptance intervals  $[-x_\rho, x_\rho]$  for normal  $N(0,1)$

$1-\rho$	0.8	0.85	0.9	0.95	0.98	0.99	0.998	0.999
$x_\rho$	1.28	1.44	1.64	1.96	2.32	2.57	3.09	3.29

➤ Thus

$$\text{Prob}\left\{-1.967 < \frac{\bar{x} - \hat{\mu}}{0.23/4} < 1.967\right\} = 0.95$$

or

$$\text{Prob}\{-0.113 < \bar{x} - \hat{\mu} < 0.113\} = 0.95$$

or

$$\text{Prob}\{1.237 < \hat{\mu} < 1.463\} = 0.95$$

- Since  $\hat{\mu} = 1.4$  lies within the above acceptance interval, **we accept**  $H_0$ , i.e.,

$$\mu = \hat{\mu} = 1.4$$

The interval  $[1.237, 1.463]$  is also known as confidence interval at the  $1-\rho=0.95$  level.

We say that: There is no **evidence** at the 5% level that the mean value is not equal to  $\hat{\mu}$

## ❖ The Unknown Variance Case

- Estimate the variance. The estimate

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

is unbiased, i.e.,

$$E[\hat{\sigma}^2] = \sigma^2$$

- Define the test statistic

$$q = \frac{\bar{x} - \mu}{\hat{\sigma} / \sqrt{N}}$$

$$\begin{aligned} \rightarrow E[\hat{\sigma}^2] &= \frac{1}{N-1} \sum_{i=1}^N E[(x_i - \bar{x})^2] \\ &= \frac{1}{N-1} \sum_{i=1}^N E[(x_i - \mu) - (\bar{x} - \mu)]^2 \\ &= \frac{1}{N-1} \sum_{i=1}^N \left( \sigma^2 + \frac{\sigma^2}{N} - 2E[(x_i - \mu)(\bar{x} - \mu)] \right) \end{aligned}$$

نمونه‌های آزمایش مستقل از هم اند:

$$\begin{aligned} E[(x_i - \mu)(\bar{x} - \mu)] &= \\ &= \frac{1}{N} E[(x_i - \mu)((x_1 - \mu) + \dots + (x_N - \mu))] \\ &= \frac{\sigma^2}{N} \end{aligned}$$

$$E[\hat{\sigma}^2] = \frac{N}{N-1} \frac{N-1}{N} \sigma^2 = \sigma^2$$

بدون بایاس

تقدیر  $q$  و تقدیر  $\sigma^2$  با درجه آزادی  $N-1$  است (بیوست A)

- This is no longer Gaussian. If  $x$  is Gaussian, then  $q$  follows a **t-distribution**, with  $N-1$  degrees of freedom

مثال 5.2 : مثال قبل را با  $\hat{\sigma}^2 = 0.23$  در نظر بگیریم.  
 $N = 16 \rightarrow N-1 = 15$   
 $p = 0.025 \rightarrow 1-p = 0.975$

- An example:

$x$  is Gaussian,  $N = 16$ , obtained from measurements,

$\bar{x} = 1.35$  and  $\hat{\sigma}^2 = (0.23)^2$ . Test the hypothesis

$$H_0 : \mu = \hat{\mu} = 1.4$$

at the significance level  $\rho = 0.025$ .

جدول بازه قبول برای تفتیش

➤ Table of acceptance intervals for t-distribution

N-1

Degrees of Freedom	1-ρ	0.9	<u>0.95</u>	<u>0.975</u>	0.99
12		1.78	2.18	2.56	3.05
13		1.77	2.16	2.53	3.01
14		1.76	2.15	2.51	2.98
<u>15</u>		1.75	2.13	<u>2.49</u>	2.95
16		1.75	2.12	2.47	2.92
17		1.74	2.11	2.46	2.90
<u>18</u>		1.73	<u>2.10</u>	2.44	2.88

➤  $\text{Prob} \left\{ -2.49 < \frac{\bar{x} - \hat{\mu}}{\hat{\sigma} / 4} < 2.49 \right\}$

$1.207 < \hat{\mu} < 1.493$

Thus,  $\hat{\mu} = 1.4$  is accepted ✓

$\frac{\bar{x} - \hat{\mu}}{\hat{\sigma} / \sqrt{N}} = \frac{1.30 - 1.4}{0.122 / 4} = -0.187$

$-2.29 < -0.187 < 2.29$

✓  $\hat{\mu} = 1.4$  ? ✓ 26

کاربرد تست  $t$  در انتخاب ویژگی :

## ❖ Application in Feature Selection

آیا ویژگی داده شده برای مسئله طبقه بندی مناسب است یا خیر؟

- The goal here is to test against **zero** the **difference**  $\mu_1 - \mu_2$  of the respective means in  $\omega_1, \omega_2$  of a single feature.

$$E[x] = \mu_1$$

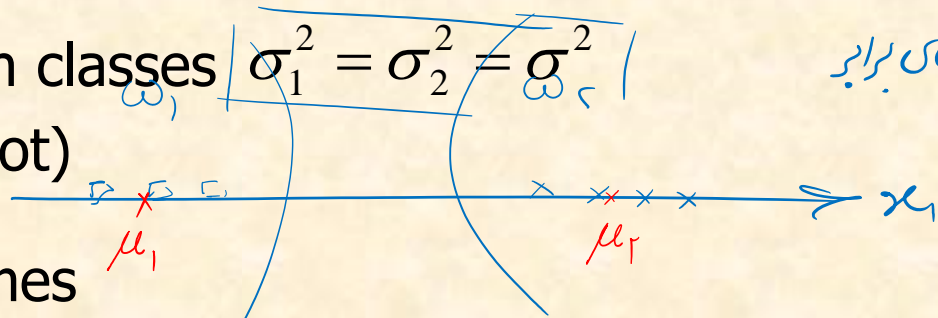
- Let  $\underline{x}_i$   $i=1, \dots, N$ , the values of a feature in  $\underline{\omega_1}$

$$E[y] = \mu_2$$

- Let  $\underline{y}_i$   $i=1, \dots, N$ , the values of the same feature in  $\underline{\omega_2}$

- Assume in both classes  $|\sigma_1^2 = \sigma_2^2 = \sigma_c^2|$  (unknown or not)

در بایس های برابر



- The test becomes

$$H_0 : \Delta\mu = \mu_1 - \mu_2 = \underline{\underline{0}}$$

$$\checkmark H_1 : \Delta\mu \neq 0$$

$$z = x - y$$

$$E[z] = \mu_1 - \mu_2$$



➤ Define

$$z = x - y$$

➤ Obviously

$$E[z] = \mu_1 - \mu_2$$

➤ Define the average

$$\bar{z} = \frac{1}{N} \sum_{i=1}^N (x_i - y_i) = \bar{x} - \bar{y}$$

➤ Known Variance Case: Define

حالت واریانس معلوم :

آماره تست :  
test statistics

$$q = \frac{(\bar{x} - \bar{y}) - (\hat{\mu}_1 - \hat{\mu}_2)}{\sigma \sqrt{\frac{2}{N}}}$$

➤ This is  $N(0,1)$  and one follows the procedure as before.

5.1



- **Unknown Variance Case:**  
Define the test statistic

حالت واریانس نامعلوم:

$$q = \frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{S_z \sqrt{\frac{2}{N}}}$$

$$S_z^2 = \frac{1}{2N - 2} \left( \sum_{i=1}^N (x_i - \bar{x})^2 + \sum_{i=1}^N (y_i - \bar{y})^2 \right) = \frac{1}{2} (\hat{\sigma}_1^2 + \hat{\sigma}_2^2)$$

- $q$  is t-distribution with  $2N-2$  degrees of freedom, 5.2
- Then apply appropriate tables as before.

مثال: مقدار  $q$  در دو جدول به صورت زیر داده شده است:

- **Example:** The values of a feature in two classes are:

$\omega_1$ : 3.5, 3.7, 3.9, 4.1, 3.4, 3.5, 4.1, 3.8, 3.6, 3.7

$\omega_2$ : 3.2, 3.6, 3.1, 3.4, 3.0, 3.4, 2.8, 3.1, 3.3, 3.6

آیا این فرق با سطح اعتماد  $\rho = 0.05$  فرق مناسبی برای طبقه بندی در جدول است؟

Test if the mean values in the two classes differ significantly, at the significance level  $\rho = 0.05$

➤ We have

$$\rightarrow \omega_1: \bar{x} = 3.73, \hat{\sigma}_1^2 = 0.0601$$

$$\rightarrow \omega_2: \bar{y} = 3.25, \hat{\sigma}_2^2 = 0.0672$$

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j; \quad \hat{\sigma}_1^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2$$

For  $N=10$

$$S_z^2 = \frac{1}{2} (\hat{\sigma}_1^2 + \hat{\sigma}_2^2) = \frac{1}{2} (0.0601 + 0.0672) = 0.06365$$

$$q = \frac{(\bar{x} - \bar{y}) - 0}{S_z \sqrt{\frac{2}{10}}}$$

$$q = 4.25$$

➤ From the table of the t-distribution with  $2N-2=18$  degrees of freedom and  $\rho=0.05$ , we obtain  $D=[-2.10, 2.10]$  and since  $q=4.25$  is outside  $D$ ,  $H_1$  is accepted and the feature is selected.

## ❖ Class Separability Measures

The emphasis so far was on individually considered features. However, such an approach cannot take into account existing correlations among the features. That is, **two features may be rich in information, but if they are highly correlated we need not consider both of them.** To this end, in order to search for possible correlations, we consider features **jointly** as elements of **vectors**. To this end:

چگونگی درجه‌ها را بصورت جداگانه در نظر گرفتیم. در این سمت توجه‌ها را بصورت توأم بررسی می‌کنیم.

صاف :

- Discard poor in information features, by means of a statistical test.

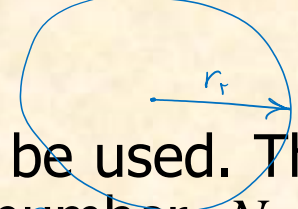
۱- حذف ویژگی‌های ضعیف با ارزیابی انفرادی آنها توسط آمار

۲- مشورت کردن نسبت به نگه‌داشتن یا حذف آن‌ها (A)

$$A_1 = \pi r_1^2$$



$$A_2 = \pi r_2^2$$



- Choose the maximum number,  $l$ , of features to be used. This is dictated by the specific problem (e.g., the number,  $N$ , of available training patterns and the type of the classifier to be adopted).

(A) تحت نظر خوب

➤ Combine remaining features to search for the "best" combination. To this end:

۳- ترکیب ویژگی‌های باقی‌مانده به صورت بردارهای لاتین  
 و جستجوی بهترین بردار

- Use different feature combinations to form the feature vector. Train the classifier, and choose the combination resulting in the best classifier performance.

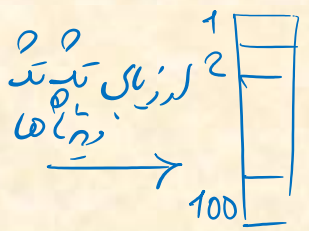
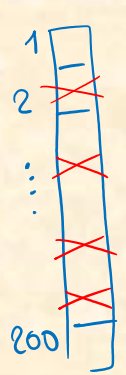
A major **disadvantage** of this approach is the high complexity. Also, local minima, **may** give misleading results.

عیب این روش پیچیده می‌باشد زیرا آن است.

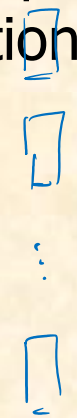
- Adopt a class separability measure and choose the best feature combination against this cost.

تعریف معیارهای جدایی بین کلاس‌ها  
 که به عنوان تابع هزینه انتخاب  
 در مرحله استفاده می‌شود.

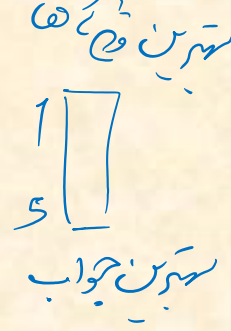
ویژگی‌های گزیده



$l=5$   
 ترکیب‌های  
 5 تایی  
 ویژگی‌ها



انتخاب بهترین  
 بردار 5 تایی



➤ **Class separability measures:** Let  $\underline{x}$  be the current feature combination vector.

5.6.1  
دیوژرانش

- **Divergence.** To see the rationale behind this cost, consider the two – class case. Obviously, if on the **average** the value of  $\ln \frac{p(\underline{x} | \omega_1)}{p(\underline{x} | \omega_2)}$  is close to zero, then  $\underline{x}$  should be a poor feature combination. Define:

$$- D_{12} = \int_{-\infty}^{+\infty} p(\underline{x} | \omega_1) \ln \frac{p(\underline{x} | \omega_1)}{p(\underline{x} | \omega_2)} d\underline{x}$$

$$- D_{21} = \int_{-\infty}^{+\infty} p(\underline{x} | \omega_2) \ln \frac{p(\underline{x} | \omega_2)}{p(\underline{x} | \omega_1)} d\underline{x}$$

$$- d_{12} = D_{12} + D_{21}$$

$$p(\omega_1 | \underline{x}) > p(\omega_2 | \underline{x})$$

$$\frac{p(\omega_1 | \underline{x})}{p(\omega_2 | \underline{x})} > 1$$

$$\ln \frac{p(\omega_1 | \underline{x})}{p(\omega_2 | \underline{x})} > 0$$

مجموعه  
البتار دیوژرانش  
در حالت دو کلاس

$d_{12}$  is known as the **divergence** and can be used as a class separability measure.

- For the multi-class case, define  $d_{ij}$  for every pair of classes  $\omega_i, \omega_j$  and the **average divergence** is defined as

1  
2  
3

$d_{12}$   
 $d_{13}$   
 $d_{23}$

دوریاتش خندونه

$$d = \sum_{i=1}^M \sum_{j=1}^M P(\omega_i)P(\omega_j)d_{ij}$$

- Some properties: **Kullback-Leibler** دوریاتش یا فاصله

$$\left\{ \begin{array}{l} d_{ij} \geq 0 \\ d_{ij} = 0, \text{ if } i = j \\ d_{ij} = d_{ji} \end{array} \right.$$

- **Large** values of  $d$  are indicative of **good** feature combination.

هرچه  $d$  بزرگتر باشد نشانه خوبی است. مناسب است.



ماتریس کوورتانس برای تقوین گاوسی:

$$\Sigma_z = \Sigma_{z|z} + \Sigma_{z|\mu}$$

$$d_{z|z} = \text{---} \quad (5.22)$$

برای حالت یقینی:

$$d_{z|z} = \frac{1}{2} \left( \frac{\sigma_z^2}{\sigma_z^2} + \frac{\sigma_z^2}{\sigma_z^2} - 2 \right) + \frac{1}{2} (\mu_z - \mu_z)^2 \left( \frac{1}{\sigma_z^2} + \frac{1}{\sigma_z^2} \right)$$

ساده شده

نتیجه: میانجی‌های بندری تنها به میانجی در دسترس است و به واریانس نیز باید بستگی داشته باشد.

$$\Sigma_z = \Sigma_{z|z} = \Sigma \text{ در حالت خاص} \quad (5.22)$$

$$d_{z|z} = (\mu_z - \mu_z)^T \Sigma^{-1} (\mu_z - \mu_z)$$

فاصله با هالانوبیس بین بردارهای میانجی متاخر و متاخرهای زودتر

نتیجه: رابطه بین کوورتانس و حتمی‌بودن رابطه مستقیم است.

$$\hat{d}_{z|z} = 2(1 - e^{-\frac{d_{z|z}}{2}})$$

کوورتانس پس از بافت

transformed

divergence

5.6.2 فاصله باتجاریه Chernoff Bound & Battacharyya  
 Distance حد شریف

کمینه احتمال خطای قابل استیابی توسط طبقه بندی بین برای دو احتمال  $\omega_i$  و  $\omega_j$ :

$$P_e = \int_{-\infty}^{+\infty} \min \left[ \underbrace{p(\omega_i) p(x|\omega_i)}_a, \underbrace{p(\omega_j) p(x|\omega_j)}_b \right] dx \quad (5.23)$$

میانگین احتمال خطای فوق در حالت کلی همین است. ولی در توفیق چه تلاقی آن را به دست آوریم.

$$\min[a, b] \leq a^s b^{1-s} \quad \text{for } a, b > 0, 0 \leq s \leq 1 \quad (5.24)$$

$$(5.23), (5.24) \rightarrow P_e \leq p(\omega_i)^s p(\omega_j)^{1-s} \int_{-\infty}^{+\infty} p(x|\omega_i)^s p(x|\omega_j)^{1-s} dx \equiv \epsilon_{CB}$$

با کمینه کردن حد شریف بستن  $s$  می توانیم به کمینه احتمال خطای دست پیدا کرد.

در حالت خاص این حد به ازای  $s = 1/2$  به صورت زیر بدست می آید:

$$P_e \leq \epsilon_{CB} = \sqrt{p(\omega_i) p(\omega_j)} \int_{-\infty}^{+\infty} \sqrt{p(x|\omega_i) p(x|\omega_j)} dx$$

در حالت  $\mathcal{P}$  می  $N(\mu_i, \Sigma_i), N(\mu_j, \Sigma_j)$

$$\downarrow \epsilon_{CB} = \sqrt{p(\omega_i) p(\omega_j)} e^{-B} \uparrow$$

Chernoff حد



فاصله  
Battacharyya

$$B = \frac{1}{\lambda} (\underline{\mu}_1 - \underline{\mu}_2)^T \left( \frac{\Sigma_1 + \Sigma_2}{2} \right)^{-1} (\underline{\mu}_1 - \underline{\mu}_2) + \frac{1}{2} \ln \frac{|\frac{\Sigma_1 + \Sigma_2}{2}|}{\sqrt{|\Sigma_1| |\Sigma_2|}}$$

۱.۱ در مین

این فاصله را در حالت به عنوان فاصله میان دو دسته‌های بزرگ می‌توان استفاده کرد.

در حالت  $\Sigma_1 = \Sigma_2 = \Sigma$  این فاصله به حد بهینه شریف دست پیدا می‌کند و فاصله باتاچاریا به پستی از فاصله ها لایوس تبدیل می‌شود.

مثال: فرض کنید  $P(\omega_1) = P(\omega_2)$  و تقویم های گاوسی  $\mathcal{N}(\underline{\mu}_1, \sigma_1^2 I)$  و  $\mathcal{N}(\underline{\mu}_2, \sigma_2^2 I)$

فاصله باتاچاریا:

$$B = \frac{1}{2} \ln \frac{(\frac{\sigma_1^2 + \sigma_2^2}{2})^l}{\sqrt{\sigma_1^{2l} \sigma_2^{2l}}} = \frac{1}{2} \ln \left( \frac{\sigma_1^2 + \sigma_2^2}{2 \sigma_1 \sigma_2} \right)^l$$

حالت خاص  $l=1$

$$\text{if } \sigma_1 = 10 \sigma_2 \rightarrow B = 0.1097$$

$$P_e \leq 0.2225$$

$$\text{if } \sigma_1 = 100 \sigma_2 \rightarrow B = 1.9561 \rightarrow P_e \leq 0.0077$$

هرچه اختلاف سوارانس ها بزرگتر باشد، حد خطای کمتر می یابید.

۹. میانگین های برابر

$$\sigma_1 = 1, \sigma_2 = 0.01$$

هرچه  $\sigma$  بزرگتر باشد حد خطای کمتر خواهد شد.

$$\frac{\sigma_2}{\sigma_1} \rightarrow 0 \rightarrow P_e \rightarrow 0$$

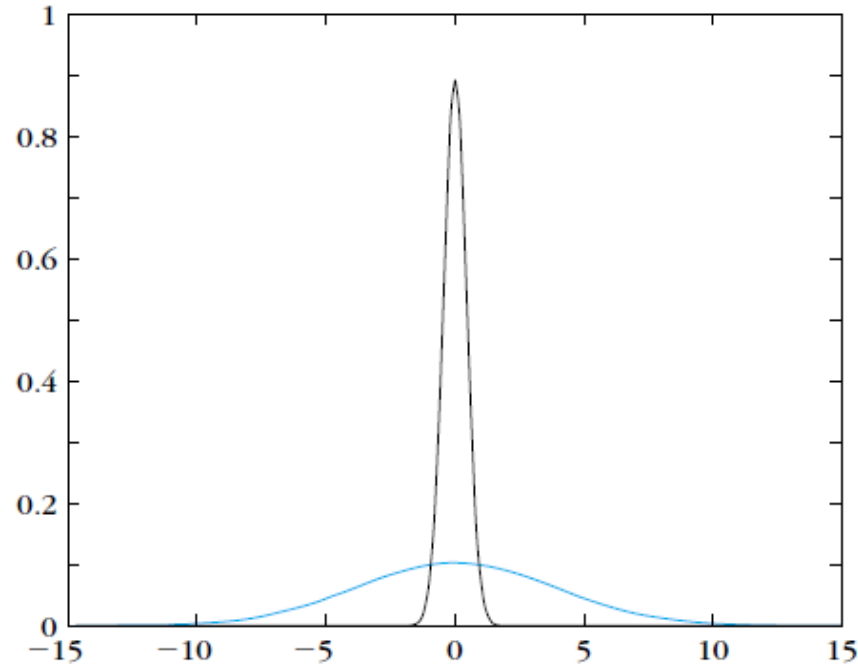


FIGURE 5.4

Gaussian pdfs with the same mean and different variances.

➤ **Scatter Matrices.** These are used as a measure of the way data are scattered in the respective feature space.

- **Within-class** scatter matrix

ماتریس پراکندگی درون کلاسی

$$S_w = \sum_{i=1}^M P_i \Sigma_i$$

چگونگی پراکندگی داده‌ها در فضای ویژگی را اندازه‌گیری می‌کند.

$\Sigma_i$ : ماتریس کوواریانس کلاسی  $\omega_i$

where

$$\Sigma_i = E \left[ \left( \underline{x} - \underline{\mu}_i \right) \left( \underline{x} - \underline{\mu}_i \right)^T \right]$$

and

$$P_i \equiv P(\omega_i) \approx \frac{n_i}{N}$$

$n_i$  the number of training samples in  $\omega_i$ .

Trace  $\{S_w\}$  is a measure of the **average variance** of the features, over all classes.

ردی ماتریس  $S_w$

$A_{n \times n}$   $tr(A) = \sum_{i=1}^n \alpha_{ii}$  مجموع عناصر قطر اصلی

- Between-class scatter matrix

ماتریس پراکندگی بین کلاس

$$S_b = \sum_{i=1}^M P_i (\underline{\mu}_i - \underline{\mu}_0) (\underline{\mu}_i - \underline{\mu}_0)^T$$

میانین میانین کلاس

$$\underline{\mu}_0 = \sum_{i=1}^M P_i \underline{\mu}_i$$

Trace  $\{S_b\}$  is a measure of the average distance of the mean of each class from the respective global one.

ردّ ط 5 معاری از فاصله میانین ، میانین هر کلاس از میانین سراسری است.

- Mixture scatter matrix

ماتریس پراکندگی مخلوط

$$S_m = E \left[ (\underline{x} - \underline{\mu}_0) (\underline{x} - \underline{\mu}_0)^T \right]$$

It turns out that:

$$\boxed{S_m = S_w + S_b}$$

➤ Measures based on Scatter Matrices.

$$\bullet \uparrow J_1 = \frac{\text{Trace}\{S_m\} \uparrow}{\text{Trace}\{S_w\} \downarrow}$$

$$\bullet J_2 = \frac{|S_m|}{|S_w|} = |S_w^{-1} S_m|$$

$$\bullet J_3 = \text{Trace}\{S_w^{-1} S_m\}$$

- Other criteria are also possible, by using various combinations of  $S_m$ ,  $S_b$ ,  $S_w$ .

The above  $J_1$ ,  $J_2$ ,  $J_3$  criteria take high values for the cases where:

- Data are clustered together within each class.
- The means of the various classes are far.

حالت خاص: یک بُعدی و دو قده سه هم احتمال

- Fisher's discriminant ratio. In one dimension and for two equiprobable classes the determinants become:

$$|S_w| \propto \sigma_1^2 + \sigma_2^2$$

$$|S_b| \propto (\mu_1 - \mu_2)^2$$

and

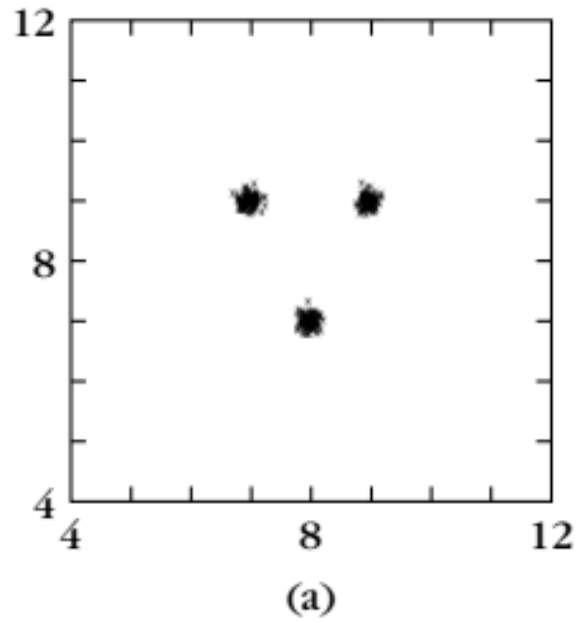
$$\frac{|S_b|}{|S_w|} = \frac{(\mu_1 - \mu_2)^2 \uparrow}{\sigma_1^2 + \sigma_2^2 \downarrow} = FDR \uparrow$$

known as **Fischer's ratio**.

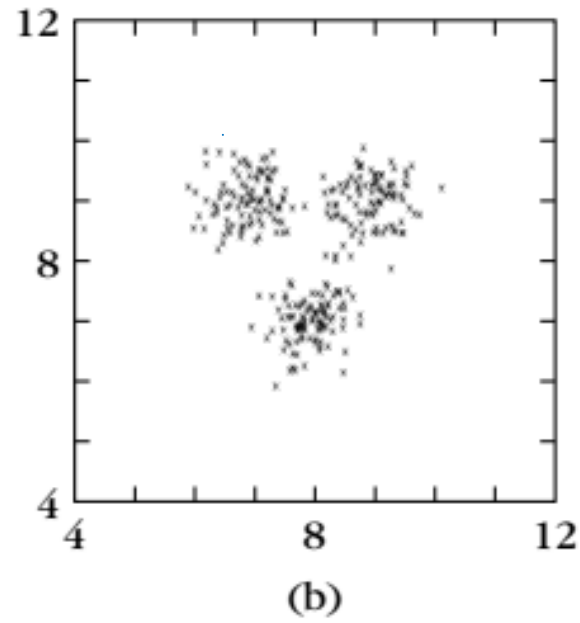
برای حالت چند قده سه

$$FDR_1 = \sum_{i=1}^M \sum_{j \neq i}^M \frac{(\mu_i - \mu_j)^2}{\sigma_i^2 + \sigma_j^2}$$

$$J_{\mu} = 194,7$$

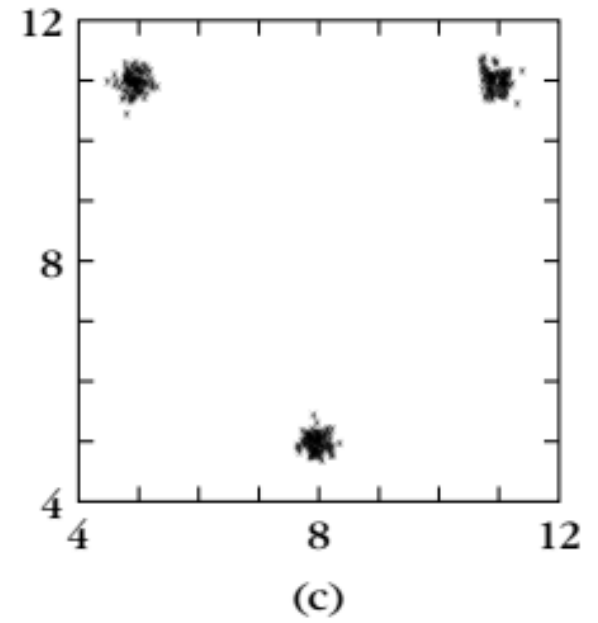


$$J_{\mu} = 112,0$$



$$\underline{\dot{J}_{\mu}}$$

$$J_{\mu} = 920,9$$





# 5.7 انتخاب زیر مجموعه و مرتب‌ها:

5.7.1: انتخاب و مرتب‌ها: تعداد بهار  $c(k)$  را می‌کند و واضح‌ها را به صورت نزولی مرتب می‌کنیم.  
 $k=1, 2, \dots, m$

ل و مرتب‌ها متناسب با بهترین ل تعداد  $c(k)$  به عنوان برادر و مرتب‌ها انتخاب می‌شوند.

$$c(k) = \min_{\text{جذب}} d_{ij}$$

حداقل دیندر این روی هر جذب است worst case

عیب: در روش انتخاب و مرتب‌ها، اکابر و مرتب‌ها به صورت جداگانه به عدد بزرگی قدردهی کنند.

مرتب: همیشه صبی‌گی کم

برای در نظر گرفتن همبستگی correlation بین و مرتب‌ها می‌توان از روش زیر استفاده کرد:

$$x_{nk} ; n=1, 2, \dots, N$$

$$k=1, 2, \dots, m$$

$k$  این و مرتب‌ها از آنوی  $n$

$$\rho_{ij} = \frac{\sum_{n=1}^N x_{ni} x_{nj}}{\sqrt{\sum_{n=1}^N x_{ni}^2 \sum_{n=1}^N x_{nj}^2}}$$

cross correlation coefficient

5.13  
 $| \rho_{ij} | \leq 1$

الگوریتم انتساب دایره با در نظر گرفتن همبستگی متقابل:

- یک مقدار جدیدی پذیرگی که مساوی  $(c)$  را در نظر می گیریم و مقدار آن را به ازای تمام دایره های موجود  $x_k$   $k=1, 2, \dots, m$

به دست می آوریم. این مقادیر به ترتیب نزدیکی مرتب می کنیم و دایره با بهترین مقدار  $c$  را انتخاب می کنیم. فرض کنید این دایره  $z_1$  باشد.

- برای انتخاب دایره  $z_2$  ضریب همبستگی متقابل که با رابطه  $z_1$  تعریف شده را بین  $z_2$  و

$(m-1)$  دایره باقیمانده می گیریم:  $z_2$  برای  $z_1 \neq z_2$

- دایره  $z_2$  را حدود انتخاب می کنیم که:

$$z_2 = \arg \max_j \{ \alpha_1 c(z_j) - \alpha_2 |z_j| \}, \text{ for all } j \neq z_1$$

$\alpha_1$  و  $\alpha_2$  ضرایب وزن دهی

به این ترتیب دایره  $z_2$  را به عنوان دایره  $z_1$  انتخاب می کنیم. علاوه بر  $c$  زیاد، همبستگی متقابل کمتری با دایره اول داشته باشد.

باینین همبستگی با دایره های قبلی

$$z_k = \arg \max_j \left\{ \alpha_1 c(z_j) - \frac{\alpha_2}{k-1} \sum_{r=1}^{k-1} |z_r| \right\} \text{ for } j \neq z_r, r=1, 2, \dots, k-1$$

- برای  $z_k, k=3, \dots, l$

متویف تابع هزینه‌ها را در هر یک از روزها تصرف می‌کند:

$$\text{Fine 83} \rightarrow \left\{ \alpha_1 c_1(z) + \alpha_2 c_2(z) - \frac{\alpha_3}{k-1} \sum_{r=1}^{k-1} |p_r z| \right\}$$

دو بهای جداگانه بزرگ

## 5.7.2 انتخاب بردار مرتبه:

- راهکار فیلتر کردن: جستجو از بین همه حالت‌ها ممکن:

$$\binom{m}{l} = \frac{m!}{l!(m-l)!} \quad \text{بهای جداگانه بزرگ}$$

مثال:  $m=20, l=5 \rightarrow \binom{20}{5} = 15,504$

مسئله پیچیده و سخت‌تر است.

- راهکار Wrapper: در این راهکار به جای بهای جداگانه بزرگ یک بهای طبقه بندی شده را در نظر می‌گیریم و ترکیبی که کمترین بهای طبقه بندی را داشته باشد را انتخاب می‌کنیم.

هر دو این روشها، شکر پیچیده تر است، در ادامه به برای حل آن مر توکل کد روش های جستجوی زیر بهینه

استفاده نمود: suboptimal Searching Techniques:

زیر بهینه

انتخاب تدریجی روبه عقب: (مکس)

sequential Backward Selection

موضوع کنید  $m=4$  ویژگی در اختیار داریم، هدف انتخاب دو ویژگی از آن هاست.  
 $x_1, x_2, x_3, x_4$

میار جدیدی بذیر کی صفت  $C$  را در نظر می گیریم و مقدمات آن را برای بردار  $x_1, x_2, x_3, x_4$  به این

صورت آوریم.  
 $[x_1, x_2, x_3, x_4]^T$ ,  $[x_1, x_2, x_3]^T$ ,  $[x_1, x_2]^T$ ,  $[x_1]^T$   
 - یک ویژگی را حذف می کنیم و چه حالتی این در نظر

میار را برای هر کدام از این چهار بردار به دست می آوریم و ترکیب بهترین مقدمات  $C$  را انتخاب می کنیم

موضوع کنید  $[x_1, x_2, x_3]^T$  انتخاب شود.  
 - یک ویژگی را حذف کرد، به حالتی این در نظر  
 $[x_1, x_2, x_3]^T$ ,  $[x_1, x_2]^T$ ,  $[x_1]^T$

میار را برای هر کدام از بهترین مقدمات را انتخاب می کنیم. مثلاً  $[x_1, x_2]^T$  بردار با دو ویژگی

در مقایسه با حالت جستجوی کامل  $(l^m)$  در این روش تعداد ترکیب‌ها از رابطه زیر به دست می‌آید:

$$1 + \frac{l}{2}(m+1)m - l(l+1) \quad \text{مثلاً } 5.15$$

- انتخاب ترتیبی رده جلو (استقیم) sequential Forward selection

- میان رده‌های  $l$  و  $l+1$  بهترین و بدترین  $l$  انتخاب می‌کنیم مثلاً  $x_1$

- برای هر ترکیب‌های دوبی که شامل  $x_1$  می‌شوند یعنی  $[x_1, x_2]^T$ ,  $[x_1, x_3]^T$ ,  $[x_1, x_4]^T$

میان رده‌ها سه گروه و ترکیب‌ها بهترین مقادیر را انتخاب می‌کنیم مثلاً  $[x_1, x_2]^T$

نکته  $l=3$  یا بیشتر باشد جدول را ادامه می‌دهیم.

در این حالت تعداد ترکیب‌ها  $l(m-l+1)/2$  می‌باشد.

سوئال: در چه حالت‌هایی از چه روش استفاده کنیم؟  
 Forward: اگر  $l$  کوچک و  $m$  بزرگ باشد  
 Backward: اگر  $l$  بزرگ و  $m$  نزدیک باشد

## ❖ Ways to combine features:

Trying to form all possible combinations of  $\ell$  features from an original set of  $m$  selected features is a computationally hard task. Thus, a number of **suboptimal** searching techniques have been derived.

- **Sequential backward selection.** Let  $x_1, x_2, x_3, x_4$  the available features ( $m=4$ ). The procedure consists of the following steps:
  - Adopt a class separability criterion (could also be the error rate of the respective classifier). Compute its value for **ALL** features considered **jointly**  $[x_1, x_2, x_3, x_4]^T$ .
  - Eliminate one feature and for each of the possible resulting combinations, that is  $[x_1, x_2, x_3]^T, [x_1, x_2, x_4]^T, [x_1, x_3, x_4]^T, [x_2, x_3, x_4]^T$ , compute the class separability criterion value  $C$ . Select the best combination, say  $[x_1, x_2, x_3]^T$ .



- From the above selected feature vector eliminate one feature and for each of the resulting combinations,  $[x_1, x_2]^T$ ,  $[x_2, x_3]^T$  compute  $[x_1, x_3]^T$  and  $C$  select the best combination.

The above selection procedure shows how one can start from  $m$  features and end up with the “best”  $\ell$  ones. Obviously, the choice is **suboptimal**. The number of required calculations is:

$$1 + \frac{1}{2}((m+1)m - \ell(\ell+1))$$

In contrast, a full search requires:

$$\binom{m}{\ell} = \frac{m!}{\ell!(m-\ell)!}$$

operations.



➤ **Sequential forward selection.** Here the reverse procedure is followed.

- Compute  $C$  for each feature. Select the “best” one, say  $x_1$
- For all possible 2D combinations of  $x_1$ , i.e.,  $[x_1, x_2]$ ,  $[x_1, x_3]$ ,  $[x_1, x_4]$  compute  $C$  and choose the best, say  $[x_1, x_3]$ .
- For all possible 3D combinations of  $[x_1, x_3]$ , e.g.,  $[x_1, x_3, x_2]$ , etc., compute  $C$  and choose the best one.

The above procedure is repeated till the “best” vector with  $\ell$  features has been formed. This is also a **suboptimal** technique, requiring:

$$\ell m - \frac{\ell(\ell - 1)}{2}$$

operations.

## - روشهای جستجو در Floating Search Methods :

روشهای قبلی در صورت حذف یک ترم یا اضافه شدن یک ترم، فرقی برای اضافه شدن یا حذف مجدد آنها وجود ندارد. به این روش اثر آت (اثر گف) معروف است  $Nesting Effect$

در این حد آن لذت بخش است در استفاده می کنیم.  $Floating Forward selection$  :

گام اول : اضافه کردن Inclusion : ویژگی نه به همراه مجموعه قبلی بهترین مقدار  $c$  را نتیجه می دهد انتخاب می کنیم }  
گام دوم : تست Test : ویژگی را با حذف شدن کمترین اثر در مجموعه قبلی ها دارد و پیدا می کنیم }  
گام سوم : حذف Exclusion : ویژگی را قبل از حذف در مرحله اولی رد می کنیم -

جزئیات بیشتر در دربارهای نذ  $P. 287$  کتاب

مقدار دهی اولیه :  $Sequential Forward$  برای  $n$  دست آوردن  $X_2$

جواب بهتر در از این همگی بیشتر

## ➤ Floating Search Methods

The above two procedures suffer from the **nesting effect**. Once a bad choice has been done, there is no way to reconsider it in the following steps.

In the floating search methods one is given the opportunity in **reconsidering a previously discarded feature or to discard a feature that was previously chosen**.

The method is still **suboptimal**, however it leads to **improved performance**, at the expense of complexity.

## ➤ Remarks:

- Besides suboptimal techniques, some optimal searching techniques can also be used, provided that the optimizing cost has certain properties, e.g., monotonic.
- Instead of using a class separability measure (**filter techniques**) or using directly the classifier (**wrapper techniques**), one can **modify** the cost function of the classifier appropriately, so that to perform feature selection and classifier design in a single step (**embedded**) method.
- For the choice of the separability measure a multiplicity of costs have been proposed, including **information theoretic costs**.

## 5.8 تولد ویژگی‌ها برینه Optimal Feature Generation

تألفون از ویژگی‌های جدیدی پذیرگی به شیوه غیرفعال استفاده شد به این معنی برای کاهش پذیرگی میزکن گام‌های دیگری  
انتخاب شده در طبقه بندی از این ویژگی‌ها بهره‌برداری.

در این بخش از ویژگی‌های جدیدی پذیرگی به شیوه فعال و در فرآیند تولد ویژگی استفاده می‌کنیم.

روش LDA  
Fisher's Linear discriminant analysis

آنالیز تمایز خطی

- حالت دو طبقه:  $x$  در فضای  $m$  بعدی از دو دسته

هدف: تولد ویژگی  $y$  به صورت ترکیب خطی از ویژگی‌های  $x$  به نحوی که ما به طبقه بندی را فشرده کرده و در فضای  
پایه کمتر حل کند.

هدف به دست آوردن جهت  $w$  در فضای  $m$  بعدی که در امتداد آن دو دسته به بهترین نحو تمیز داده می‌شوند.

$$y = w_1 x_1 + w_2 x_2 + \dots$$

$$y = \frac{w^T x}{\|w\|}$$

تصویر  $x$  در امتداد  $w$

از  $\|w\|$  به دلیل اینکه در جهت تأثیر منفی در تصویر کنیم.

Fisher  
Discriminant  
Ratio

$$FDR = \frac{(\mu_1 - \mu_2)^T}{\sigma_1^2 + \sigma_2^2}$$

$\mu_1, \mu_2, \sigma_1^2, \sigma_2^2$  به ترتیب میانگین ها و واریانس های  $\varphi$  در دسته های  $\omega_1, \omega_2$  باشند.

$$y = \underline{w}^T \underline{x} \rightarrow \mu_i = \underline{w}^T \underline{\mu}_i, \quad i=1,2$$

$$|(\mu_1 - \mu_2)^T| = \underline{w}^T (\underline{\mu}_1 - \underline{\mu}_2) (\underline{\mu}_1 - \underline{\mu}_2)^T \underline{w} \propto |\underline{w}^T S_b \underline{w}|$$

$$\sigma_i^2 = E[(y - \mu_i)^2] = E[\underline{w}^T (\underline{x} - \underline{\mu}_i) (\underline{x} - \underline{\mu}_i)^T \underline{w}] = \underline{w}^T \sum_{i=1,2} \underline{w}$$

$$|\sigma_1^2 + \sigma_2^2| \propto \underline{w}^T S_w \underline{w}$$

$$FDR = \frac{\underline{w}^T S_b \underline{w}}{\underline{w}^T S_w \underline{w}}$$

با بیشینه سازی FDR نسبت به  $\underline{w}$  به جواب می‌رسیم. این رابطه وقتی بیشینه می‌شود:

$$S_b \underline{w} = \lambda S_w \underline{w}$$

← بزرگترین مقدار ویژه  $\lambda$  ماتریس  $S_b^{-1} S_w^{-1}$

در این مسئله به ما می‌دهد:







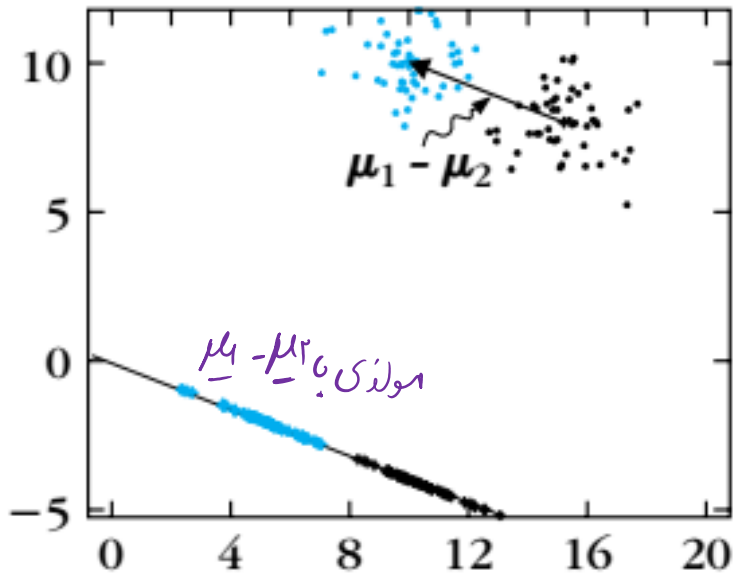
# LDA: Linear Discriminant Analysis

## 2 Classes

### 2 Feature to 1 Feature

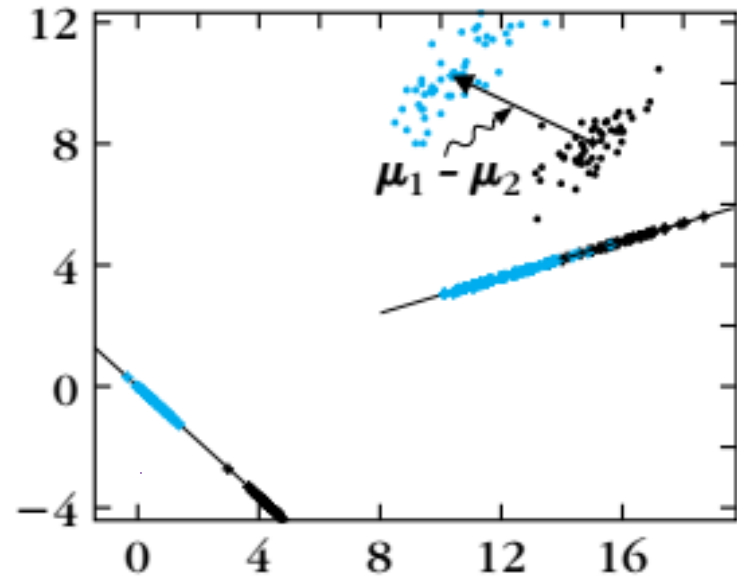
5.6

حالت کوردیناسیون قطری



(a)

حالت کوردیناسیون غیرقطری



(b)

حالت چندگانه:

$$y = A^T x$$

$\begin{matrix} \xrightarrow{m \times 1} \\ \downarrow \\ \begin{matrix} l \times 1 & l \times m \end{matrix} \end{matrix}$

$\leftarrow$   $l \times 1$        $\rightarrow$   $l \times m$

•  $J_p = \text{trace}\{S_w^{-1} S_b\}$  ، این را می بینیم

$$S_{yw} = A^T S_{xw} A \quad S_{yb} = A^T S_{xb} A$$

$$J_p(A) = \text{trace}\{(A^T S_{xw} A)^{-1} (A^T S_{xb} A)\} \quad \text{problem 5.17}$$

$$\frac{\partial J_p(A)}{\partial A} = 0 \quad \rightarrow \quad (S_{xw}^{-1} S_{xb}) A = A (S_{yw}^{-1} S_{yb})$$

$$B^T S_{yw} B = I \quad ; \quad B^T S_{yb} B = D \quad l \times l$$

مقداری

$$\hat{y} = B^T y = B^T A^T x$$

$$J_p(\hat{y}) = J_p(y)$$

(5.44)

$$(S_{xw}^{-1} S_{xb}) C = CD$$

$\begin{matrix} \downarrow \\ AB \end{matrix}$

$$\hat{y} = C^T \underline{x} \longrightarrow J_{r_{\max}} \quad : \quad l = M-1 \text{ حالت}$$

$$J_{r, \hat{y}} = J_{r, \underline{x}}$$

$$\hat{y} = (\underline{\mu}_1 - \underline{\mu}_2)^T S_{\text{XW}}^{-1} \underline{x} \quad : \quad \underline{M} = 2 \text{ حالت}$$

↓  
LDA (است)

$$J_{r, \hat{y}} < J_{r, \underline{x}} \quad : \quad l < M-1 \text{ حالت}$$

بیان هندسی:

5.7:

