# Introduction to Machine Learning  CS182

Lu Sun

School of Information Science and Technology

ShanghaiTech University
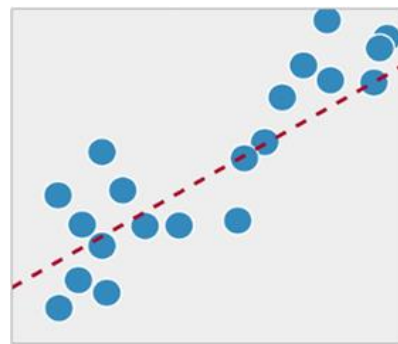
October 19, 2023

Today:
- Linear Methods for Classification I
  - Introduction
  - Linear regression of an indicator matrix
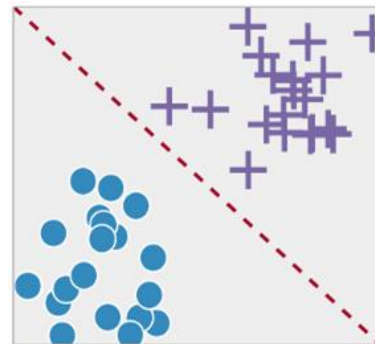  - Linear discriminant analysis

Readings:
- The Elements of Statistical Learning (ESL), Chapters 4.1, 4.2 and 4.3

# Linear Methods for Classification I

- Introduction
- Linear regression of an indicator matrix
- Linear discriminant analysis

Regression

Classification

# Introduction

## Example
Handwritten digits recognition

Input variables
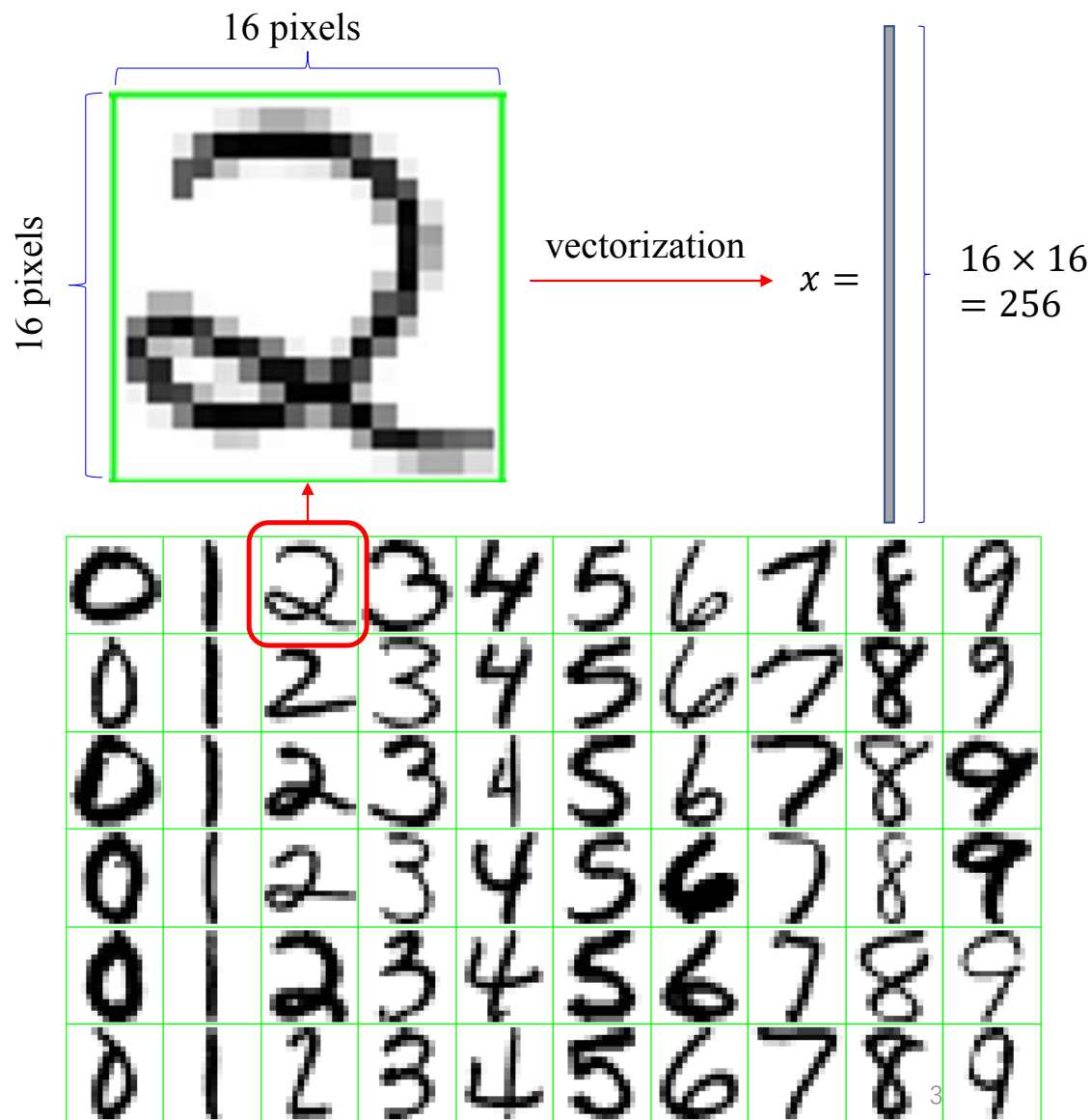$$X = (X_0, X_1, X_2, \ldots, X_{256})^T$$
Categorical output variable $G$ with values from
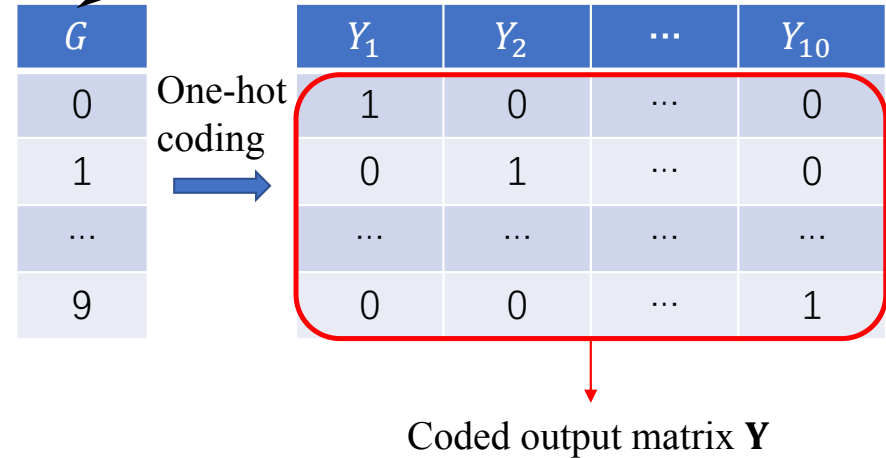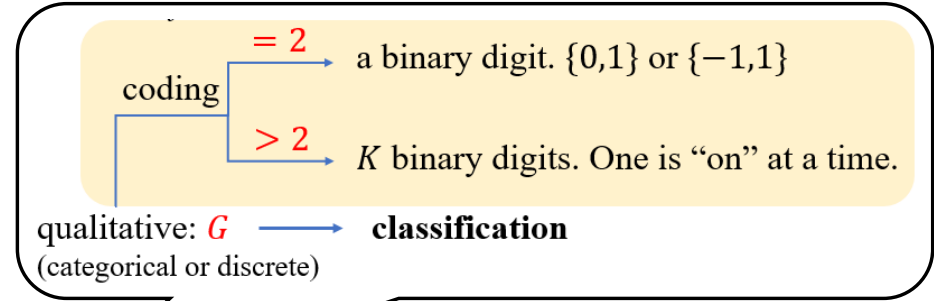$$\mathcal{G} = \{0, 1, 2 \ldots, 9\}$$

Non-binary (multi-class) classification

16 pixels

16 pixels

vectorization $\longrightarrow$ $x =$

$16 \times 16$
$= 256$



3

# Introduction

Example

Handwritten digits recognition

| | = 2 | a binary digit. {0,1} or {−1,1} |
| coding | > 2 | $K$ binary digits. One is "on" at a time. |

qualitative: $G$ $\longrightarrow$ **classification**
(categorical or discrete)

| | $X_0$ | $X_1$ | $X_2$ | ... | $X_{256}$ |
|---|---|---|---|---|---|
| sample $x_1^T$ | 1 | 0.156 | 0.432 | ... | 0.824 |
| | 1 | 0.671 | 0.014 | ... | 0.969 |
| | ... | ... | ... | ... | ... |
| | 1 | 0.523 | 0.142 | ... | 0.718 |

All-one vector for the intercept

Data matrix **X**

| $G$ |
|---|
| 0 |
| 1 |
| ... |
| 9 |

One-hot coding $\Longrightarrow$

| $Y_1$ | $Y_2$ | ... | $Y_{10}$ |
|---|---|---|---|
| 1 | 0 | ... | 0 |
| 0 | 1 | ... | 0 |
| ... | ... | ... | ... |
| 0 | 0 | ... | 1 |

Coded output matrix **Y**

$$\min_{\mathbf{B}} \|\mathbf{Y} - \mathbf{XB}\|_F^2 \quad \Longrightarrow \quad \widehat{\mathbf{B}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$$

1. Any problems?
2. Other methods?

# Introduction

Binary classification
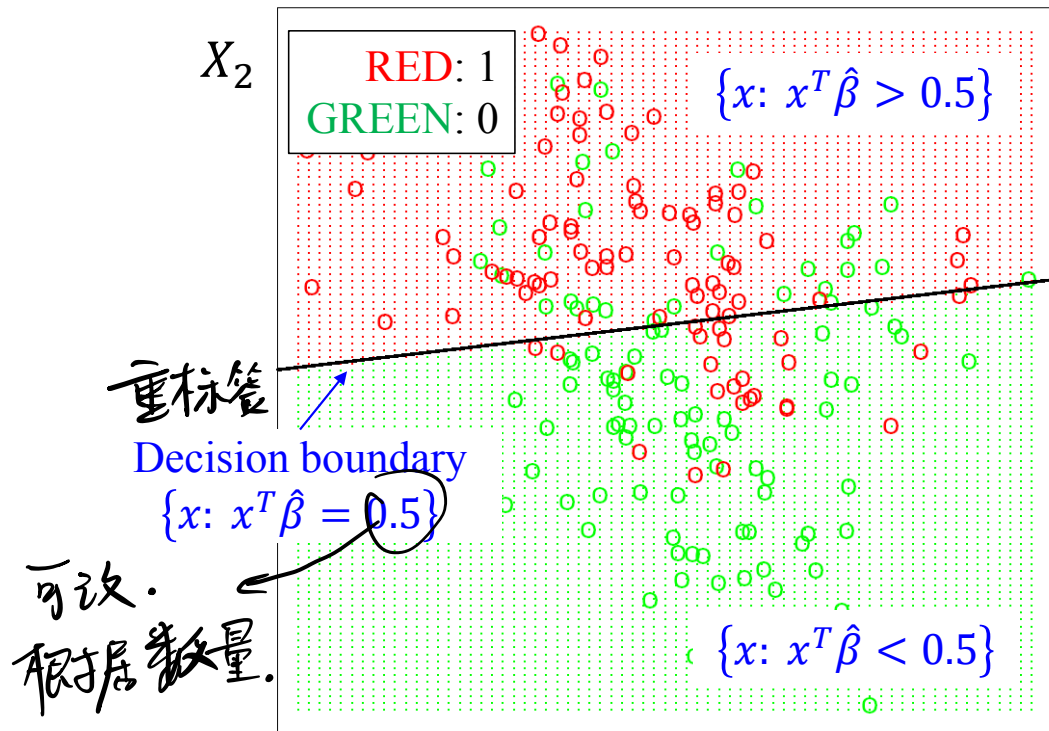
- Linear regression

$$f(x) = \beta_0 + x^T\beta$$

- Least squares solution

$$\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

- Decision boundary

$$\{x : x^T\hat{\beta} = threshold\}$$

  ❑ $threshold = 0$, if $y \in \{-1, 1\}$
  ❑ $threshold = 0.5$, if $y \in \{0, 1\}$

$X_2$

| RED: 1 |
| GREEN: 0 |

$\{x: x^T\hat{\beta} > 0.5\}$

重标签

Decision boundary
$\{x: x^T\hat{\beta} = 0.5\}$

可改.
根据数量.

$\{x: x^T\hat{\beta} < 0.5\}$

$X_1$

# Introduction

Multi-class classification

- Linear regressions for $K$ classes
$$f_k(x) = \beta_{k0} + x^T \beta_k, \qquad k = 1, \ldots, K$$

有很多种情况（$n^2$级）

- Decision boundary between classes $k$ and $\ell$:
$$\{x: \hat{f}_k(x) = \hat{f}_\ell(x)\}$$

For $K$ classes, there are $\binom{K}{2} = \frac{K(K-1)}{2}$ decision boundaries

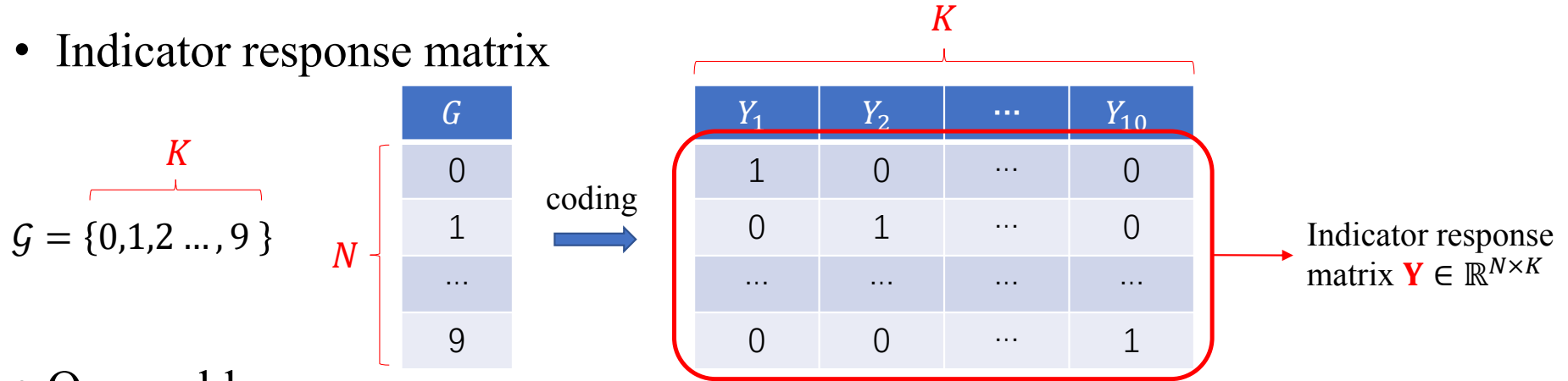- That is an affine set or hyperplane:
$$\{x: (\hat{\beta}_{k0} - \hat{\beta}_{\ell 0}) + x^T(\hat{\beta}_k - \hat{\beta}_\ell) = 0\}$$

# Linear Methods for Classification I

- Introduction
- Linear regression of an indicator matrix
- Linear discriminant analysis

# Linear Regression of an Indicator Matrix

- Indicator response matrix

$$\mathcal{G} = \{0, 1, 2 \dots, 9\}$$

| $G$ |
|:---:|
| 0 |
| 1 |
| ... |
| 9 |

coding →

| $Y_1$ | $Y_2$ | ... | $Y_{10}$ |
|:---:|:---:|:---:|:---:|
| 1 | 0 | ... | 0 |
| 0 | 1 | ... | 0 |
| ... | ... | ... | ... |
| 0 | 0 | ... | 1 |

Indicator response matrix $\mathbf{Y} \in \mathbb{R}^{N \times K}$

- Our problem:

$$\widehat{\mathbf{B}} = \underset{\mathbf{B}}{\mathrm{argmin}} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2$$

$$\mathbf{B} = (\beta_1, \beta_2, \dots, \beta_{10}) \in \mathbb{R}^{(p+1) \times K}$$

- The fitted values on $\mathbf{X}$:

$$\widehat{\mathbf{Y}} = \mathbf{X}\widehat{\mathbf{B}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

# Linear Regression of an Indicator Matrix

A new observation $x$ is classified by

- Compute the fitted output

$$\hat{f}(x) = \widehat{\mathbf{B}}^T \begin{pmatrix} 1 \\ x \end{pmatrix} = \begin{pmatrix} \hat{f}_1(x) \\ \hat{f}_2(x) \\ \vdots \\ \hat{f}_K(x) \end{pmatrix} \in \mathbb{R}^K$$

向量 => 对每个类别输出
=> 找最大值对应的类别

- Classify $x$ according to

$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \, \hat{f}_k(x)$$

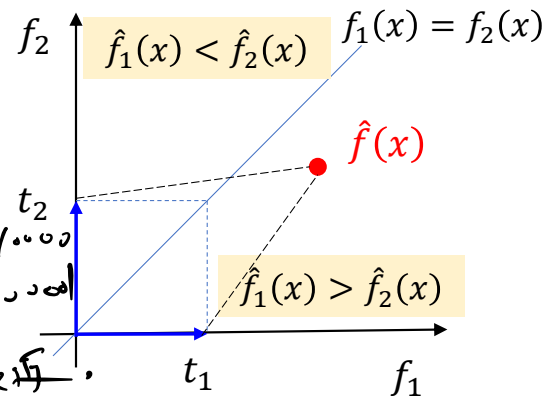独火编码如 $\begin{cases} t_1 = 1 \cdots 0 \\ \vdots \\ t_k = 0 \cdots 1 \end{cases}$

- Or equivalently,

$$\hat{G}(x) = \operatorname{argmin}_{k \in \mathcal{G}} \left\| \hat{f}(x) - t_k \right\|_2^2$$ 找最近.

where $t_k = (0, \dots, 0, 1, 0, \dots, 0)^T \in \mathbb{R}^K$ is a target with 1 being the $k$-th element

$f_2$    $\hat{f}_1(x) < \hat{f}_2(x)$     $f_1(x) = f_2(x)$

$\hat{f}(x)$

$t_2$

$\hat{f}_1(x) > \hat{f}_2(x)$

$t_1$     $f_1$

# Linear Regression of an Indicator Matrix

Categorical output variable $G$ with values from $\mathcal{G} = \{1, \dots, K\}$.

- The zero-one loss function

$$L(k, \ell) = \begin{cases} 1, & k \neq \ell \\ 0, & k = \ell \end{cases}$$

- Expected prediction error (EPE) w.r.t. $\Pr(G, X)$

$$\text{EPE} = \text{E}\left[ L\left( G, \hat{G}(X) \right) \right]$$

- Pointwise minimization leads to

$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\text{argmin}} \sum_{\ell=1}^{K} L(k, \ell) \Pr(G = \ell | X = x)$$

$$= \underset{k \in \mathcal{G}}{\text{argmax}} \; \boxed{\Pr(G = k | X = x)} \longleftarrow \text{posterior}$$

# Linear Regression of an Indicator Matrix

A new observation $x$ is classified by

- Compute the fitted output

$$\hat{f}(x) = \hat{B}^T \begin{pmatrix} 1 \\ x \end{pmatrix} = \begin{pmatrix} \hat{f}_1(x) \\ \hat{f}_2(x) \\ \vdots \\ \hat{f}_K(x) \end{pmatrix} \in \mathbb{R}^K$$
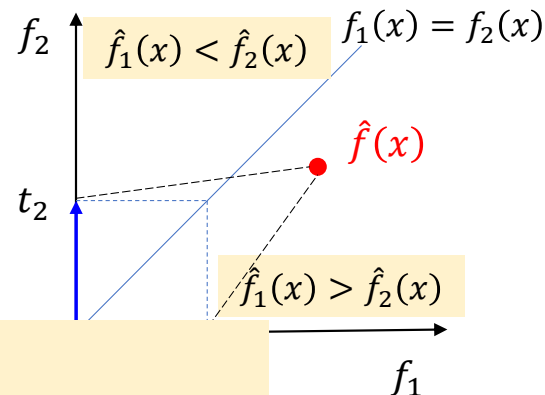
- Classify $x$ according to

$$\boxed{\hat{G}(x) = \underset{k \in \mathcal{G}}{\arg\max}\, \hat{f}_k(x)}$$

- Minimizing EPE w.r.t. the 0-1 loss gives rise to

$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\arg\max}\, \Pr(G = k | X = x)$$

- Our question:

    Are the $\hat{f}_k(x)$ reasonable estimates of the posterior $\Pr(G = k | X = x)$?



$f_2$    $\hat{f}_1(x) < \hat{f}_2(x)$    $f_1(x) = f_2(x)$

$\hat{f}(x)$

$t_2$

$\hat{f}_1(x) > \hat{f}_2(x)$

$f_1$

ment

# Linear Regression of an Indicator Matrix

?

Linear classification:
$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \hat{f}_k(x)$$

Minimizing EPE:
$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \Pr(G = k | X = x)$$

Two defining properties of probability
1. $\sum P = 1$
2. $0 < P < 1$

- It can be verified that $\sum_{k \in \mathcal{G}} \hat{f}_k(x) = 1$
- However, it is possible that $\hat{f}_k(x) < 0$ or $\hat{f}_k(x) > 1$

可取负或 +10

span (col (X))

将截矩吸入

Suppose that $\mathbf{X} \leftarrow (\mathbf{1}_N, \mathbf{X})$ and
$$\widehat{\mathbf{Y}} = \hat{f}(\mathbf{X}) = \mathbf{X}\widehat{\mathbf{B}} = \left( \hat{f}_1(\mathbf{X}), \dots, \hat{f}_K(\mathbf{X}) \right)$$

We have the followings     维度为k的列向量
$$\sum_{k=1}^K \hat{f}_K(\mathbf{X}) = \widehat{\mathbf{Y}} \cdot \mathbf{1}_K$$

Indicator matrix
标准矩阵
$$= \mathbf{X}\widehat{\mathbf{B}} \cdot \mathbf{1}_K$$
$$= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} \cdot \mathbf{1}_K$$
$$= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \cdot \mathbf{1}_N$$
$$= \mathbf{H} \cdot \mathbf{1}_N$$

$\mathbf{H} \cdot \mathbf{1}_N$ is a projection of $\mathbf{1}_N$ onto the
column space of $\mathbf{X}$, thus $\mathbf{H} \cdot \mathbf{1}_N = \mathbf{1}_N$

# Linear Regression of an Indicator Matrix

?

Linear classification:
$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \hat{f}_k(x)$$

Minimizing EPE:
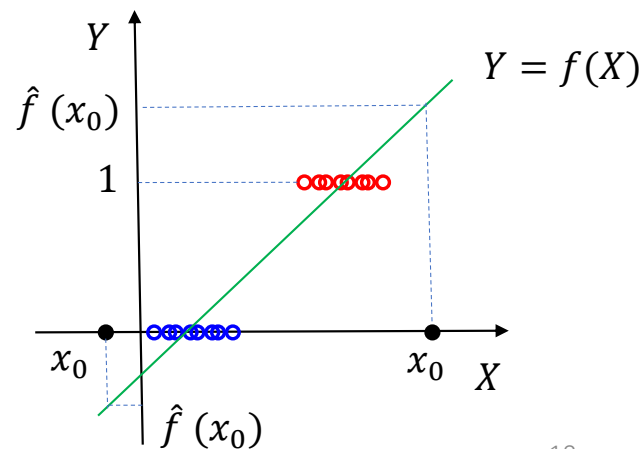$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \Pr(G = k | X = x)$$

Two defining properties of probability
1. $\sum P = 1$
2. $0 < P < 1$

- It can be verified that $\sum_{k \in \mathcal{G}} \hat{f}_k(x) = 1$
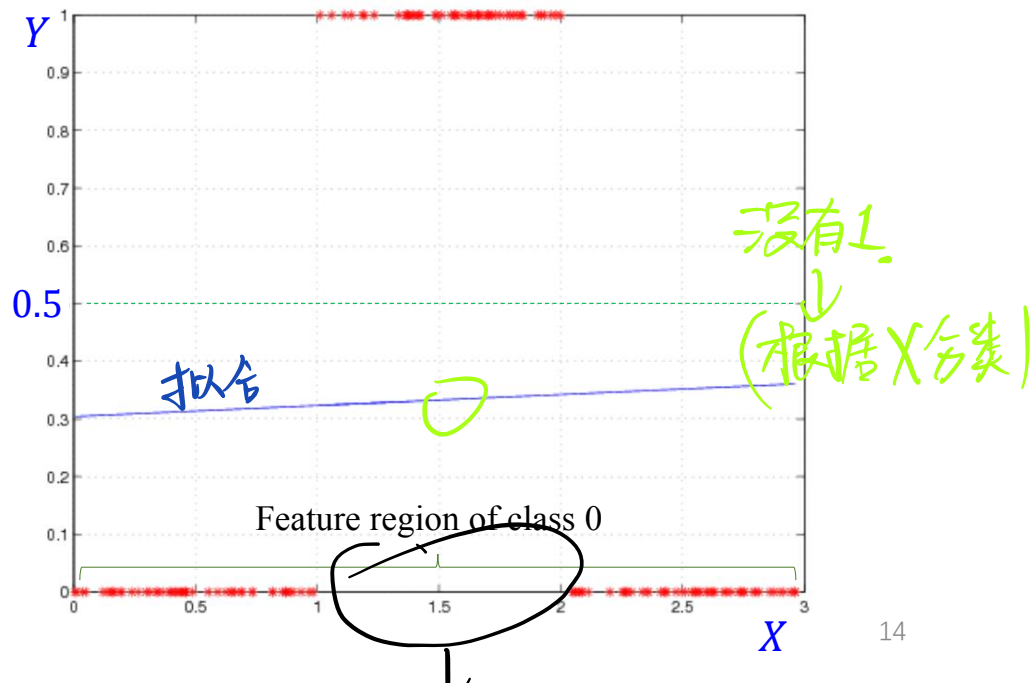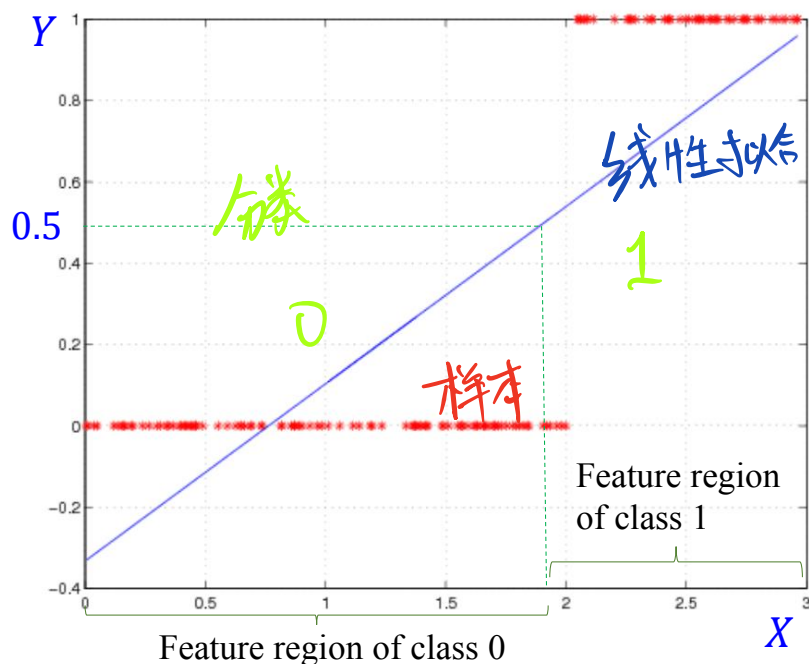- However, it is possible that $\hat{f}_k(x) < 0$ or $\hat{f}_k(x) > 1$

It possibly suffers from the problem of masking
- a class may be masked by others, i.e., there is no region in the feature space that is labeled as this class



13

# The Phenomenon of Masking

- A class may be masked by others, i.e., there is no region in the feature space that is labeled as this class
- The linear regression model is too rigid



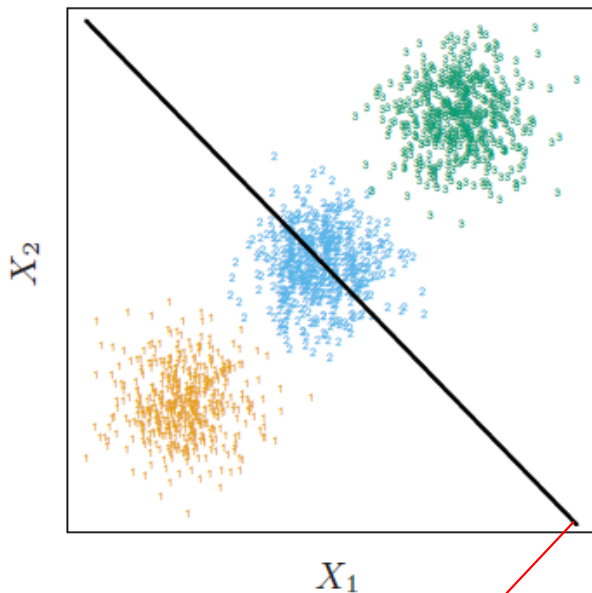Feature region of class 0

Feature region of class 1

Feature region of class 0

# The Phenomenon of Masking

这里会被遮掉.
变成全为 0 的行集

- 3-class classification



**Linear Regression**

**Linear Discriminant Analysis** ← Ideal result

Yellow: class 1
Blue: class 2
Green: class 3

$X_2$

$X_1$

$X_2$

$X_1$

Decision boundary between classes 2 and 3

Decision boundary between classes 1 and 2

The decision boundaries between 1 and 2 and between 2 and 3 are the same, so we would never predict class 2.

15

# The Phenomenon of Masking

- **3-class classification**

找哪个值大 就是哪一类
⇒ 不信有蓝色类.

The indicator matrix
$$g = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \rightarrow \mathbf{Y} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$
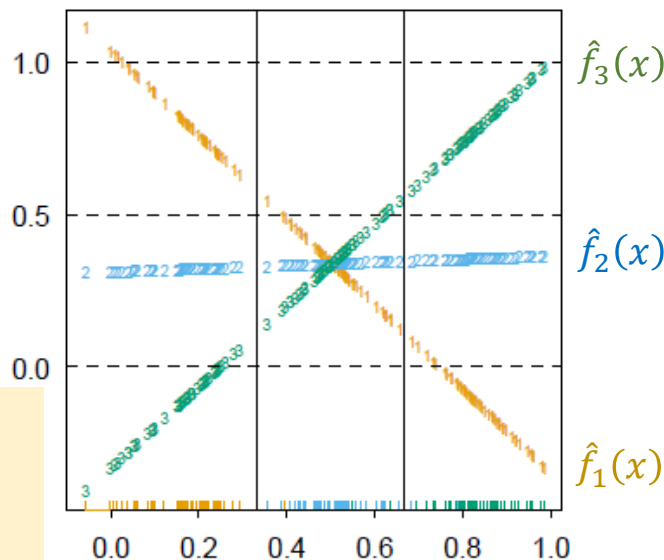
**Yellow**: class 1
**Blue**: class 2
**Green**: class 3

$$\widehat{\mathbf{B}} = \underset{\mathbf{B}}{\text{argmin}} \|\mathbf{Y} - \mathbf{XB}\|_F^2,$$
where $\mathbf{X} = (\mathbf{1}_N, \mathbf{x})$

$$\hat{f}(x) = \widehat{\mathbf{B}}^T \begin{pmatrix} 1 \\ x \end{pmatrix} = \begin{pmatrix} \hat{f}_1(x) \\ \hat{f}_2(x) \\ \hat{f}_3(x) \end{pmatrix}$$



Degree = 1; Error = 0.33

Linear regression
$Y = \beta_0 + \beta X$

Degree = 2; Error = 0.04

拓展到非线性 Quadratic regression
$Y = \beta_0 + \beta_1 X + \beta_2 X^2$

# Linear Methods for Classification I

- Introduction
- Linear regression of an indicator matrix
- **Linear discriminant analysis**

# Linear Discriminant Analysis

- Recall our discussion on linear regression of an indicator matrix

Linear classification:
$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \; \hat{f}_k(x)$$

Minimizing EPE:
$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \; \Pr(G = k | X = x)$$

- It is inappropriate to represent a posterior directly by a linear function.

# Linear Discriminant Analysis

- Idea:

  model the posterior $\Pr(G = k|X = x)$ based on the Bayes theorem

- Posterior

  用 gauss分布

  $$\Pr(G = k|X = x) = \frac{\Pr(X=x|G=k)\Pr(G=k)}{\Pr(X=x)} = \frac{\boxed{\Pr(X=x|G=k)}\boxed{\Pr(G=k)}}{\sum_{\ell=1}^{K} \Pr(X=x|G=\ell)\Pr(G=\ell)}$$

  是常数，则非

  □ Density of $X$ in class $G = k$:
  $$f_k(x) = \Pr(X = x|G = k) \quad \begin{cases} 1 & X=k \\ 0 & X\neq k \end{cases}$$

  $\mathbb{1}_{x=k} = \begin{cases} 1 & X=k \\ 0 & X\neq k \end{cases}$

  □ Class prior:

  $$\pi_k = \Pr(G = k)$$

  $$\Pr(G = k|X = x) = \frac{f_k(x)\pi_k}{\sum_{\ell=1}^{K} f_\ell(x)\pi_\ell}$$

  伯努利 Bernoulli: $\pi^x(1-\pi)^{1-x}$   多取值 Categorical: $\prod_{k=1}^{} \pi_k^{\mathbb{1}_{x=k}}$   线性    样本

- It produces LDA, QDA (quadratic DA), MDA (mixture DA), kernel DA and naïve Bayes, under various assumptions on $f_k(x)$

  判断 $\{x|\Pr(G=k|X=x) = \Pr(G=\ell|X=x)\}$ => $\frac{\Pr(k|\text{万})}{} = 1$ <=> $\frac{\Pr(\text{万}|k)\,p(k)}{}$ =1

# Linear Discriminant Analysis

$P \in (0,1) \to LDA \in (-\infty, +\infty)$

$$Pr(G = k|X = x) = \frac{f_k(x)\pi_k}{\sum_{\ell=1}^{K} f_\ell(x)\pi_\ell}$$

- Assumptions in LDA
    1. Model each class density as multivariate Gaussian

    高维 Gaussian  （了维向量）

    $$f_k(x) = \frac{1}{(2\pi)^{p/2}|\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1}(x - \mu_k)\right)$$

    2. Assume that classes share a common covariance $\Sigma_k = \Sigma, \forall k$

    假设: 任意 k, $\Sigma_k$ 相互相等.

- Compare two classes $k$ and $\ell$

Logit:

$$\log \frac{Pr(G = k|X = x)}{Pr(G = \ell|X = x)} = \log \frac{f_k(x)}{f_\ell(x)} + \log \frac{\pi_k}{\pi_\ell}$$

（了代入 $f_k(x)$）

$$= \log \frac{\pi_k}{\pi_\ell} - \frac{1}{2}(\mu_k + \mu_\ell)^T \Sigma^{-1}(\mu_k - \mu_\ell)$$
$$+ x^T \Sigma^{-1}(\mu_k - \mu_\ell),$$

Quadratic term
vanished due to
the common
covariance

只有一个 $x$, 是关于 $x$ 的线性

Decision boundary is linear w.r.t. $X$

20

顶有二次项，但因为有被消掉抵消了

# Linear Discriminant Analysis

- Parameter estimation

$\hat{\pi}_k = N_k/N$, where $N_k$ is the number of class-$k$ observations;

$\hat{\mu}_k = \sum_{g_i=k} x_i/N_k$;

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

Pooled covariance (合并方差)

$$\hat{\Sigma} = \frac{(N_1 - 1)\hat{\Sigma}_1 + (N_2 - 1)\hat{\Sigma}_2 + \cdots + (N_K - 1)\hat{\Sigma}_K}{(N_1 - 1) + (N_2 - 1) + \cdots + (N_K - 1)}, \text{where } \hat{\Sigma}_k = \frac{\sum_{g_i=k}(x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T}{N_k - 1}$$

Weighted average

不是 $N_k$，保证无偏估计．

# Linear Discriminant Analysis

Data    Class

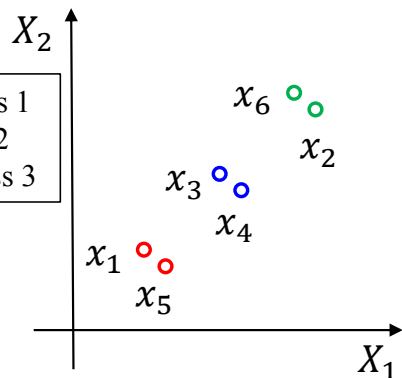| | $X_1$ | $X_2$ | $G$ |
|---|---|---|---|
| $x_1^T$ | 0.2 | 0.3 | 1 |
| $x_2^T$ | 0.8 | 0.7 | 3 |
| $x_3^T$ | 0.4 | 0.6 | 2 |
| $x_4^T$ | 0.6 | 0.4 | 2 |
| $x_5^T$ | 0.3 | 0.2 | 1 |
| $x_6^T$ | 0.7 | 0.8 | 3 |

- Class prior

$$\hat{\pi}_1 = \hat{\pi}_2 = \hat{\pi}_3 = \frac{1}{3} = \frac{N_k}{N}$$

- Class-specific sample mean

$\Gamma = 维空间$
$= ]$ $= 维向量$

$$\hat{\mu}_1 = \frac{1}{2}(x_1 + x_5) = \frac{1}{2}\begin{pmatrix}0.2\\0.3\end{pmatrix} + \frac{1}{2}\begin{pmatrix}0.3\\0.2\end{pmatrix} = \begin{pmatrix}0.25\\0.25\end{pmatrix}$$

$$\hat{\mu}_2 = \frac{1}{2}(x_3 + x_4) = \frac{1}{2}\begin{pmatrix}0.4\\0.6\end{pmatrix} + \frac{1}{2}\begin{pmatrix}0.6\\0.4\end{pmatrix} = \begin{pmatrix}0.5\\0.5\end{pmatrix}$$

$$\hat{\mu}_3 = \frac{1}{2}(x_2 + x_6) = \frac{1}{2}\begin{pmatrix}0.8\\0.7\end{pmatrix} + \frac{1}{2}\begin{pmatrix}0.7\\0.8\end{pmatrix} = \begin{pmatrix}0.75\\0.75\end{pmatrix}$$

- Common covariance

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

$$\frac{\begin{pmatrix}0.005 & -0.005\\-0.005 & 0.005\end{pmatrix} + \begin{pmatrix}0.02 & -0.02\\-0.02 & 0.02\end{pmatrix} + \begin{pmatrix}0.005 & -0.005\\-0.005 & 0.005\end{pmatrix}}{6-3} = \begin{pmatrix}0.03 & -0.03\\-0.03 & 0.03\end{pmatrix}$$

$X_2$

Green: class 1
Blue: class 2
Yellow: class 3

$x_6$
$x_2$
$x_3$
$x_4$
$x_1$
$x_5$

$X_1$

# Linear Discriminant Analysis

**Data**

| | $X_1$ | $X_2$ |
|---|---|---|
| $x_1^T$ | 0.2 | 0.3 |
| $x_2^T$ | 0.8 | 0.7 |
| $x_3^T$ | 0.4 | 0.6 |
| $x_4^T$ | 0.6 | 0.4 |
| $x_5^T$ | 0.3 | 0.2 |
| $x_6^T$ | 0.7 | 0.8 |

**Class**

| $G$ |
|---|
| 1 |
| 3 |
| 2 |
| 2 |
| 1 |
| 3 |

- For classes 1 and 2

$$\hat{\Sigma}_\lambda = \hat{\Sigma} + \lambda \mathbf{I} \longleftarrow \lambda = 1$$
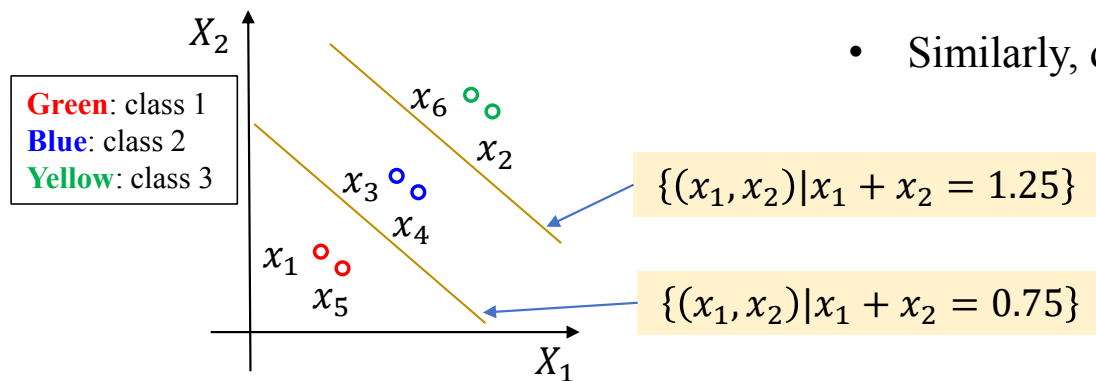
$$\log \frac{\Pr(G=1|X=x)}{\Pr(G=2|X=x)}$$

$$= \log \frac{\hat{\pi}_1}{\hat{\pi}_2} - \frac{1}{2}(\hat{\mu}_1 + \hat{\mu}_2)^T \hat{\Sigma}_\lambda^{-1}(\hat{\mu}_1 - \hat{\mu}_2) + x^T \hat{\Sigma}_\lambda^{-1}(\hat{\mu}_1 - \hat{\mu}_2)$$

$$= \frac{1}{2}(0.75, 0.75)\begin{pmatrix} 0.972 & 0.028 \\ 0.028 & 0.972 \end{pmatrix}\begin{pmatrix} 0.25 \\ 0.25 \end{pmatrix} - (x_1, x_2)\begin{pmatrix} 0.972 & 0.028 \\ 0.028 & 0.972 \end{pmatrix}\begin{pmatrix} 0.25 \\ 0.25 \end{pmatrix}$$

$$= 0.1875 - (x_1, x_2)\begin{pmatrix} 0.25 \\ 0.25 \end{pmatrix} = 0$$

- Decision boundary 1-2: $\{(x_1, x_2)|x_1 + x_2 = 0.75\}$

- Similarly, decision boundary 2-3: $\{(x_1, x_2)|x_1 + x_2 = 1.25\}$

**Green**: class 1
**Blue**: class 2
**Yellow**: class 3

$\{(x_1, x_2)|x_1 + x_2 = 1.25\}$

$\{(x_1, x_2)|x_1 + x_2 = 0.75\}$

# Linear Discriminant Analysis

- Suppose that $\log \dfrac{\Pr(G=k|X=x)}{\Pr(G=\ell|X=x)} = \delta_k(x) - \delta_\ell(x)$

  - $\delta_k(x) > \delta_\ell(x)$, class k
  - $\delta_k(x) < \delta_\ell(x)$, class $\ell$
  - $\delta_k(x) = \delta_\ell(x)$, decision boundary

线性判别函数：算每个类别记 $\delta_k$ 看哪个最大.

- Linear discriminant functions

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

$\nearrow$ Constant

Classify to class $k$ that maximizes the discriminant function

$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \, \delta_k(x)$$

Any difference?

Linear classification:
$$\hat{G}(x) = \underset{k \in \mathcal{G}}{\operatorname{argmax}} \, \hat{f}_k(x)$$

# Linear Discriminant Analysis

- Binary classification ($K = 2$)
    - ☐ Correspondence between LDA and linear classification

- Multi-class classification ($K \geq 3$)
    - ☐ LDA is different with linear classification
    - ☐ Avoid the masking problem



Degree = 1; Error = 0.33

Yellow: class 1
Blue: class 2
Green: class 3

Class 2 is masked by classes 1 and 3

# Linear Discriminant Analysis
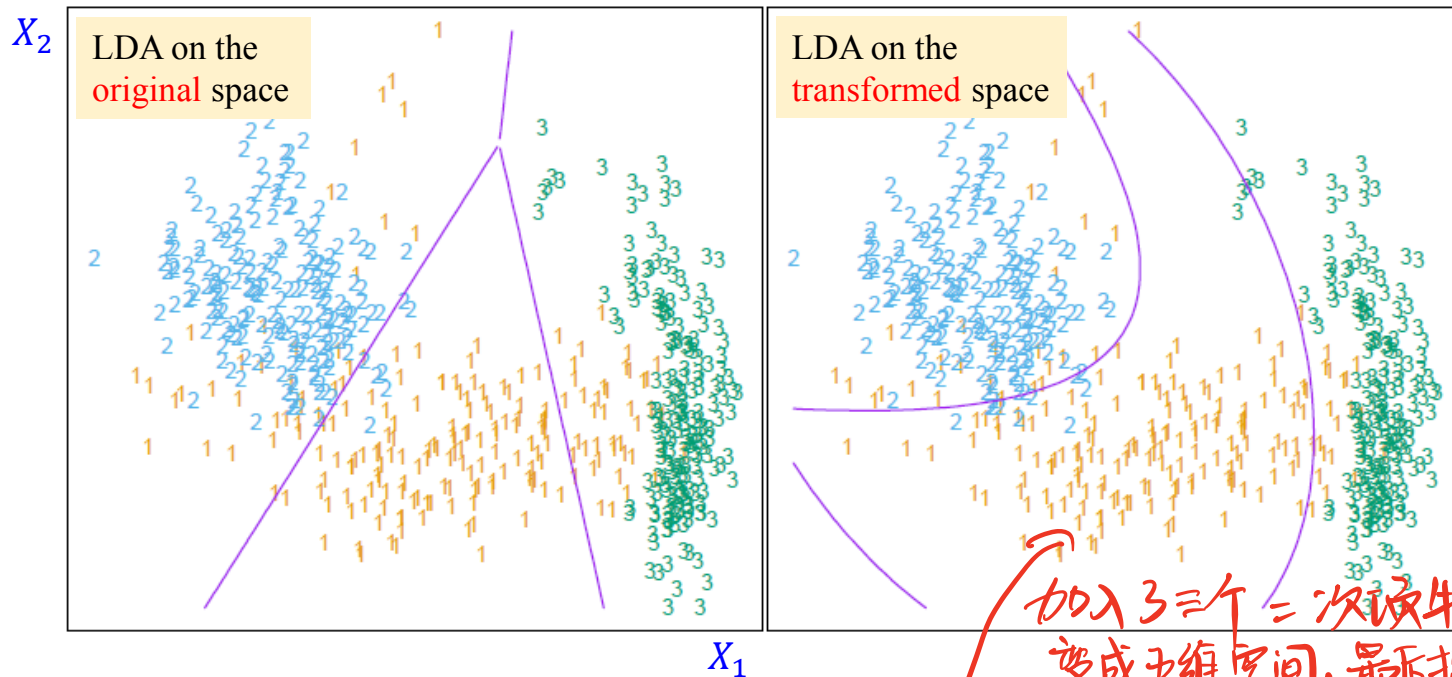


FIGURE 4.1. *The left plot shows some data from three classes, with linear decision boundaries found by linear discriminant analysis. The right plot shows quadratic decision boundaries. These were obtained by finding linear boundaries in the five-dimensional space* $X_1, X_2, X_1 X_2, X_1^2, X_2^2$. *Linear inequalities in this space are quadratic inequalities in the original space.*

取哪些二次项是一个超参数.

# Linear Discriminant Analysis



**FIGURE 4.5.** *The left panel shows three Gaussian distributions, with the same covariance and different means. Included are the contours of constant density enclosing 95% of the probability in each case. The Bayes decision boundaries between each pair of classes are shown (broken straight lines), and the Bayes decision boundaries separating all three classes are the thicker solid lines (a subset of the former). On the right we see a sample of 30 drawn from each Gaussian distribution, and the fitted LDA decision boundaries.*

# Quadratic Discriminant Analysis

特征表达更好但需要计算的更多. (相对 LDA)

- **Assumption**: Each class has a specific covariance $\Sigma_k$

- Quadratic discriminant functions

$$\delta_k(x) = -\frac{1}{2}\log|\mathbf{\Sigma}_k| - \frac{1}{2}(x - \mu_k)^T \mathbf{\Sigma}_k^{-1}(x - \mu_k) + \log \pi_k.$$

- The quadratic decision boundary between two classes $k$ and $\ell$

$$\{x : \delta_k(x) = \delta_\ell(x)\}$$

- **Difference with LDA**
  - $\Sigma_k$ has to be estimated for each class
  - LDA need to estimate $K \times p$ + $p \times p$ parameters
  - QDA need to estimate $K \times p$ + $K \times p \times p$ parameters

$\mu_k, k = 1, \dots, K$

$\Sigma$

需估计 $\pi, \mu, \Sigma$

$\Sigma_k, k = 1, \dots, K$

需要 K 个 $\Sigma$

28

# Quadratic Discriminant Analysis



**FIGURE 4.6.** *Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision boundaries for the data in Figure 4.1 (obtained using LDA in the five-dimensional space* $X_1, X_2, X_1 X_2, X_1^2, X_2^2$*). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.*

# Summary

- Linear regression of an indicator matrix
  - The indicator matrix
  - Prediction is conducted by $\hat{G}(x) = \text{argmax}_k \hat{f}_k(x)$
  - Suffer from the masking problem

- Linear discriminant analysis
  - Logit transformation: $\text{logit}(\text{Pr}(x)) = \log\left(\frac{\text{Pr}(x)}{1-\text{Pr}(x)}\right)$
  - Model the posterior $\text{Pr}(G = k|X = x)$
  - Assumptions on $\text{Pr}(X = x|G = k)$
  - Discriminant functions $\delta_k(x)$

- Quadratic discriminant analysis
  - Difference with LDA

# Classification

Simple and straightforward

Theoretical

**Linear regression**

$\min_f \text{EPE}$

Squared error loss

Zero-one loss

$\mathcal{G} = \{1, 2 \ldots, K\}$

**Indicator matrix** $\mathbf{Y} = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$

**Regression function**

$f(x) = \text{E}(Y|X = x)$

**Bayes classifier**

$\hat{G}(x) = \underset{k \in \mathcal{G}}{\text{argmax}}\, \Pr(G = k | X = x)$

Multi-output regression

Linear

Nonlinear

$(0,1) \rightarrow (-\infty, +\infty)$

**Prediction**

$\hat{f}(x) = \hat{\mathbf{B}}^T \begin{pmatrix} 1 \\ x \end{pmatrix} = \begin{pmatrix} \hat{f}_1(x) \\ \hat{f}_2(x) \\ \vdots \\ \hat{f}_K(x) \end{pmatrix}$

$\hat{G}(x) = \underset{k \in \mathcal{G}}{\text{argmax}}\, \hat{f}_k(x)$

**Least squares**

**Nearest neighbors**

**Logit transformation** $\text{logit}(x) = \log\left(\frac{x}{1-x}\right)$

**Regression**

Pairwise odds = 1

Limitation

**The masking problem** $(K \geq 3)$

**Decision boundary** $\log\dfrac{\Pr(G = k|X = x)}{\Pr(G = \ell|X = x)} = 0$

Bayes theorem

Linear boundary

**LDA, QDA, RDA**

**Logistic regression**