

Introduction to Machine Learning CS182

Lu Sun

School of Information Science and Technology

ShanghaiTech University

October 24, 2023

Today:

- Linear Methods for Classification II
 - Generalization of LDA
 - Logistic Regression → 家庭是分类
 - Summary

Readings:

- The Elements of Statistical Learning (ESL), Chapters 4.3, 4.4, 18.1, 18.2 and 18.3

Linear Methods for Classification II

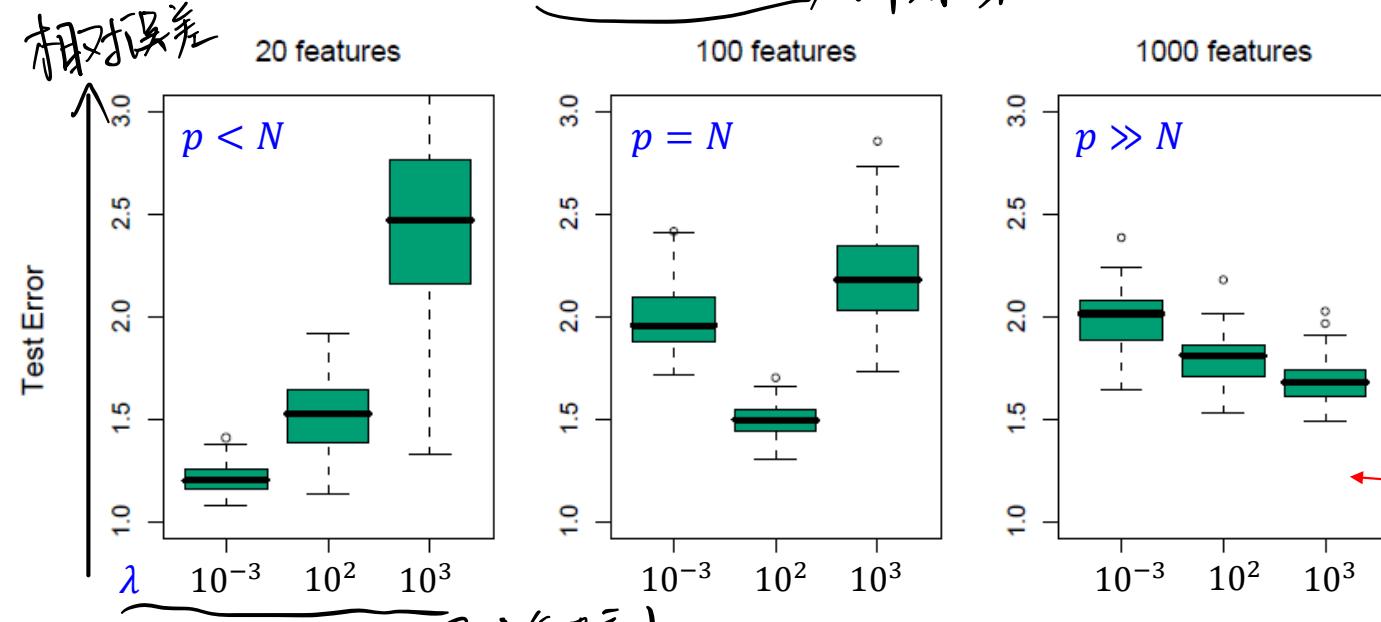
- Generalization of LDA
 - Regularized Discriminant Analysis 正則化
 - Fisher's Formulation of Discriminant Analysis
 - Logistic Regression
 - Summary
- Fisher/尼氏判斷法.*

Regularized Discriminant Analysis

高维少数据
奇异矩阵.

High dimensional problems ($p \gg N$)

- genomics problem, signal/image analysis
- Less fitting is better
 \rightarrow 训练集小, 防止过拟合.



Example

- 100 samples are generated by a linear model
- Ridge regression
- Relative error (divide by Bayes error)

维度越高模型越简单越好

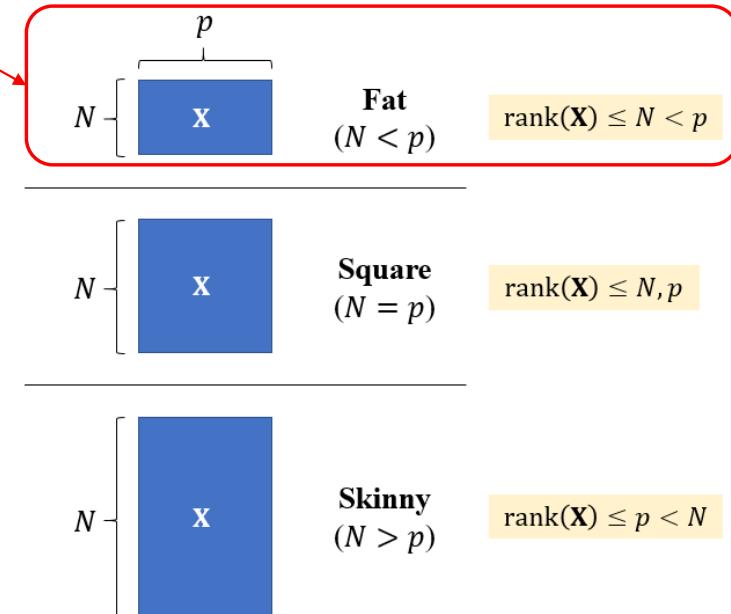
No enough information to estimate the high-dimensional covariance matrix

↑↑ 德国↑↑ 复杂↓

Regularized Discriminant Analysis

High dimensional problems ($p \gg N$)

- Cannot fit LDA to the data
 - inversion of a $p \times p$ covariance matrix Σ
 - Σ is singular, due to $\text{rank}(\Sigma) \leq N \ll p$
- Regularization is necessary
 - No enough data to estimate feature dependencies
 - E.g., independent assumption on features
 - Diagonal within-class covariance matrix
 - #paras: $K \times p \times p \rightarrow K \times p$



训练集过大时可假设协方差矩阵是对角矩阵
↓
↓

Regularized Discriminant Analysis

Regularized LDA (RLDA)

- Shrinks $\hat{\Sigma}$ towards its diagonal

$$\hat{\Sigma}(\gamma) = \gamma\hat{\Sigma} + (1 - \gamma)\text{diag}(\hat{\Sigma}), \gamma \in [0, 1]$$

where $\text{diag}(\hat{\Sigma})$ denotes a diagonal matrix sharing the same diagonal elements with $\hat{\Sigma}$

Diagonal LDA

- Independent assumption on feature dependencies

$$\hat{\Sigma} = \text{diag}(\hat{\Sigma})$$

Regularized Discriminant Analysis

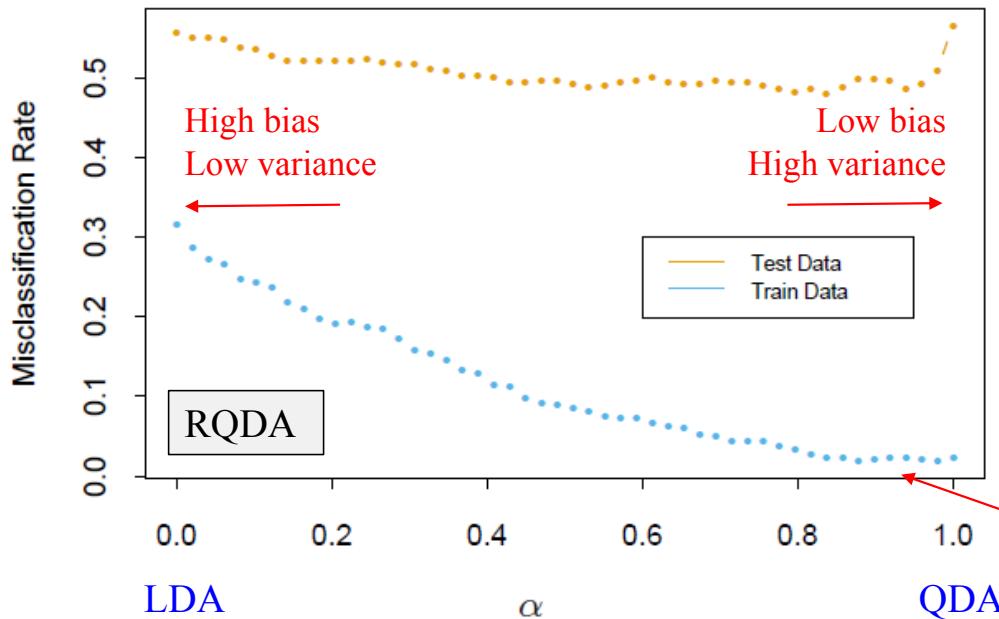
A brief summary of generalized LDA ($\alpha, \gamma \in [0, 1]$)

	Method	Covariance matrix	Effect
Linear	Regularized LDA (RLDA)	$\widehat{\Sigma}(\gamma) = \gamma\widehat{\Sigma} + (1 - \gamma)\text{diag}(\widehat{\Sigma})$	Shrink $\widehat{\Sigma}$ towards $\text{diag}(\widehat{\Sigma})$
	Diagonal LDA	$\widehat{\Sigma} = \text{diag}(\widehat{\Sigma})$	Make features independent $\Leftrightarrow \frac{1}{\lambda_i} \propto 0$
Quadratic	Regularized QDA (RQDA)	$\widehat{\Sigma}_k(\alpha) = \alpha\widehat{\Sigma}_k + (1 - \alpha)\widehat{\Sigma}$ $\Leftrightarrow \text{认为 } \Sigma_k \text{ 相同} \Leftrightarrow \lambda \propto 1$	Shrink $\widehat{\Sigma}_k$ towards $\widehat{\Sigma}$ (LDA + QDA)
	Variant of RQDA	$\widehat{\Sigma}_k(\alpha, \gamma) = \alpha\widehat{\Sigma}_k + (1 - \alpha)\widehat{\Sigma}(\gamma)$	Shrink $\widehat{\Sigma}_k$ towards $\widehat{\Sigma}(\gamma)$ (RLDA + QDA)

Regularized Discriminant Analysis

Regularized Discriminant Analysis on the Vowel Data

<https://hastie.su.domains/ElemStatLearn/>



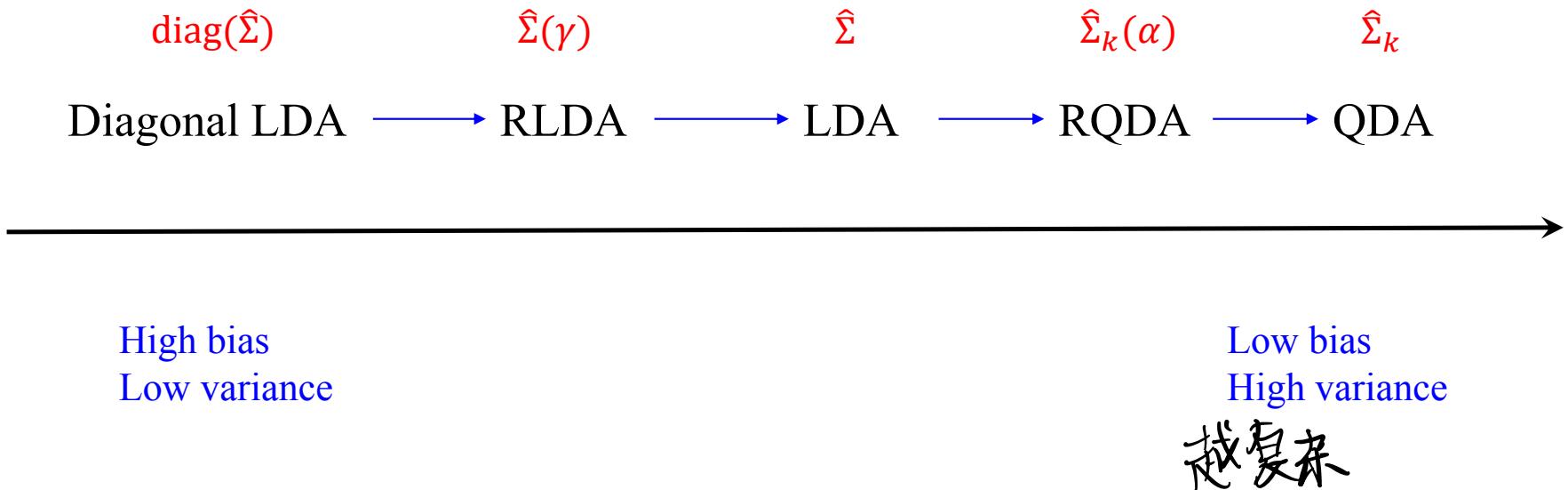
RQDA:
 $\hat{\Sigma}_k(\alpha) = \alpha\hat{\Sigma}_k + (1 - \alpha)\hat{\Sigma}$

- $\alpha = 0$, LDA
- $\alpha = 1$, QDA

The optimal model
A compromise between
QDA and LDA

FIGURE 4.7. Test and training errors for the vowel data, using regularized discriminant analysis with a series of values of $\alpha \in [0, 1]$. The optimum for the test data occurs around $\alpha = 0.9$, close to quadratic discriminant analysis.

Regularized Discriminant Analysis



计算 LDA

要看处于哪个类别内。
只考虑距圆心距离即可

Fisher's Formulation of Discriminant Analysis

LDA: Approach 1

1. Estimating $\hat{\Sigma}$, $\hat{\mu}_k$ and $\hat{\pi}_k$

2. Discriminant function

$$\delta_k(x) = x^T \hat{\Sigma}^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \hat{\Sigma}^{-1} \hat{\mu}_k + \log \hat{\pi}_k$$

3. Classify to class k that maximizes the discriminant function

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

协方差矩阵

$$\hat{\Sigma} = \frac{(X - \bar{\mu})(X - \bar{\mu})^T}{n-1}$$

$$\begin{aligned} & \text{求} \quad \hat{\Sigma} = \frac{\sum_i (\bar{x}_i - \bar{\mu})(\bar{x}_i - \bar{\mu})^T}{n-1} \\ & = \frac{\bar{x}(\sum_i x_i - \sum_i \bar{x})}{n-1} (\bar{x}^T - \bar{x}^T \bar{\mu})^T \end{aligned}$$

LDA: Approach 2

1. Estimating $\hat{\Sigma}$, $\hat{\mu}_k$ and $\hat{\pi}_k$

2. Eigen-decomposition:

$$\hat{\Sigma} = \mathbf{U} \mathbf{D} \mathbf{U}^T \quad \left\{ \begin{array}{l} \mathbf{U}: \text{正交矩阵} \\ \mathbf{D}: \text{对角} \end{array} \right.$$

3. Data spherling ($\hat{\Sigma}^* = \mathbf{I}$) 直化 (球化)

$$x^* = \mathbf{D}^{-\frac{1}{2}} \mathbf{U}^T x = \hat{\Sigma}^{-\frac{1}{2}} x \quad \text{假设协方差矩阵变成 I}$$

$$\hat{\mu}_k^* = \mathbf{D}^{-\frac{1}{2}} \mathbf{U}^T \hat{\mu}_k = \hat{\Sigma}^{-\frac{1}{2}} \hat{\mu}_k$$

4. Classify to its closest class centroid in the transformed space

$$\hat{G}(x) = \operatorname{argmin}_{k \in \mathcal{G}} \frac{1}{2} \|x^* - \hat{\mu}_k^*\|^2 - \ln \hat{\pi}_k$$

协方差

协方差

$$= \frac{\sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T}{n-1} = I$$

复杂度降低

Fisher's Formulation of Discriminant Analysis

$$1. \log \frac{\Pr(G=k|X=x)}{\Pr(G=\ell|X=x)} = \delta_k(x) - \delta_\ell(x)$$

$$2. \delta_k(x) \propto \log \Pr(G = k|X = x)$$

$$\begin{aligned} 3. \log \Pr(G = k|X = x) &= -\frac{1}{2}(x - \hat{\mu}_k)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_k) + \log \hat{\pi}_k + C \\ &= -\frac{1}{2}(x - \hat{\mu}_k)^T \mathbf{U} \mathbf{D}^{-\frac{1}{2}} (\mathbf{U} \mathbf{D}^{-\frac{1}{2}})^T (x - \hat{\mu}_k) + \log \hat{\pi}_k + C \\ &= -\frac{1}{2} \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{U}^T x - \mathbf{D}^{-\frac{1}{2}} \mathbf{U}^T \hat{\mu}_k \right)^T \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{U}^T x - \mathbf{D}^{-\frac{1}{2}} \mathbf{U}^T \hat{\mu}_k \right) + \log \hat{\pi}_k + C \\ &= -\frac{1}{2} (x^* - \hat{\mu}_k^*)^T (x^* - \hat{\mu}_k^*) + \log \hat{\pi}_k + C \\ &= -\frac{1}{2} \|x^* - \hat{\mu}_k^*\|^2 + \ln \hat{\pi}_k + C \end{aligned}$$

如果每种类别样本数相同

$$4. \hat{G}(x) = \operatorname{argmax}_{k \in \mathcal{G}} \delta_k(x) = \operatorname{argmax}_{k \in \mathcal{G}} \log \Pr(G = k|X = x) = \operatorname{argmin}_{k \in \mathcal{G}} \frac{1}{2} \|x^* - \hat{\mu}_k^*\|^2 - \ln \hat{\pi}_k$$

因为是常数 \Rightarrow 不用管。

$$\mathcal{N}(\hat{\mu}_k, \hat{\Sigma})$$

$$\hat{\pi}_k$$

$$\Pr(G = k|X = x) = \frac{\Pr(X = x|G = k)\Pr(G = k)}{\Pr(X = x)}$$

则无需修正

Fisher's Formulation of Discriminant Analysis

LDA: Approach 1

1. Estimating $\widehat{\Sigma}$, $\widehat{\mu}_k$ and $\widehat{\pi}_k$

2. Discriminant function

$$\delta_k(x) = x^T \widehat{\Sigma}^{-1} \widehat{\mu}_k - \frac{1}{2} \widehat{\mu}_k^T \widehat{\Sigma}^{-1} \widehat{\mu}_k + \log \widehat{\pi}_k$$

3. Classify to class k that maximizes the discriminant function

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

Complexity
 $\mathcal{O}(p^3)$

- Two approaches have almost the same time and storage complexity
- Approach 2 shows the potential of LDA for dimension reduction

LDA: Approach 2

1. Estimating $\widehat{\Sigma}$, $\widehat{\mu}_k$ and $\widehat{\pi}_k$

2. Eigen-decomposition:

可以在*这里* $\widehat{\Sigma} = \mathbf{U} \mathbf{D} \mathbf{U}^T$
保留一些维度

3. Data spherling ($\widehat{\Sigma}^* = \mathbf{I}$)

- $x^* = \mathbf{D}^{-\frac{1}{2}} \mathbf{U}^T x = \widehat{\Sigma}^{-\frac{1}{2}} x$
- $\widehat{\mu}_k^* = \mathbf{D}^{-\frac{1}{2}} \mathbf{U}^T \widehat{\mu}_k = \widehat{\Sigma}^{-\frac{1}{2}} \widehat{\mu}_k$

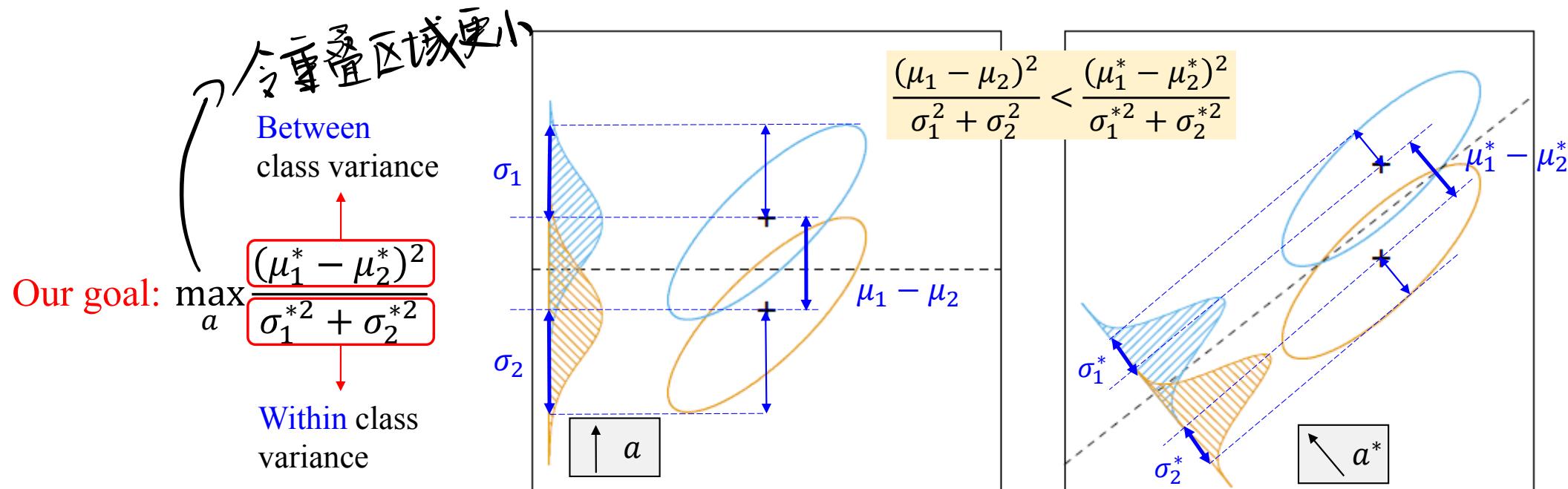
4. Classify to its closest class centroid in the transformed space

$$\hat{G}(x) = \operatorname{argmin}_{k \in \mathcal{G}} \frac{1}{2} \|x^* - \widehat{\mu}_k^*\|^2 - \ln \widehat{\pi}_k$$

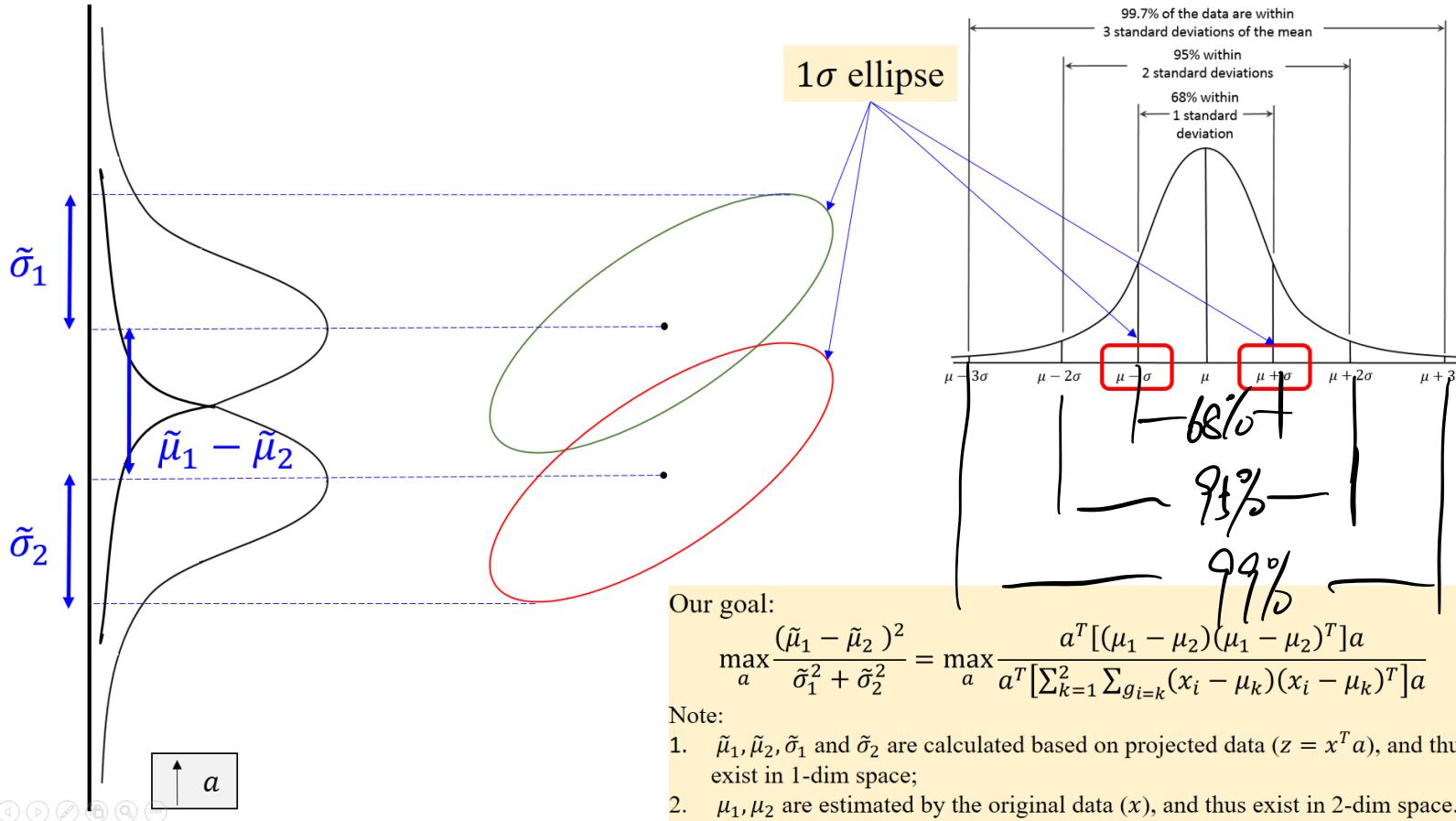
1. 两个类别二类因降维后的维度数不可超过 k-1

Fisher's Formulation of Discriminant Analysis

- Find $z = x^T a$ such that the **between class** variance is maximized relative to the **within class** variance.



Fisher's Formulation of Discriminant Analysis



$$\tilde{z}_i = \tilde{\mu}^T \alpha. \quad \begin{matrix} \tilde{z}_i = \tilde{\mu}^T \alpha \\ \tilde{\mu} \in \mathbb{R}^P \end{matrix} \quad \begin{matrix} M \in \mathbb{R}^{P \times P} \\ \Sigma \in \mathbb{R}^{P \times P} \end{matrix} \quad \begin{matrix} \hat{M} = M^T \alpha \\ \hat{\sigma}^2 = \alpha^T \Sigma \alpha \end{matrix}$$

Fisher's Formulation of Discriminant Analysis

- Maximize the Rayleigh quotient:

$$\max_a \frac{a^T \mathbf{B} a}{a^T \mathbf{W} a}$$

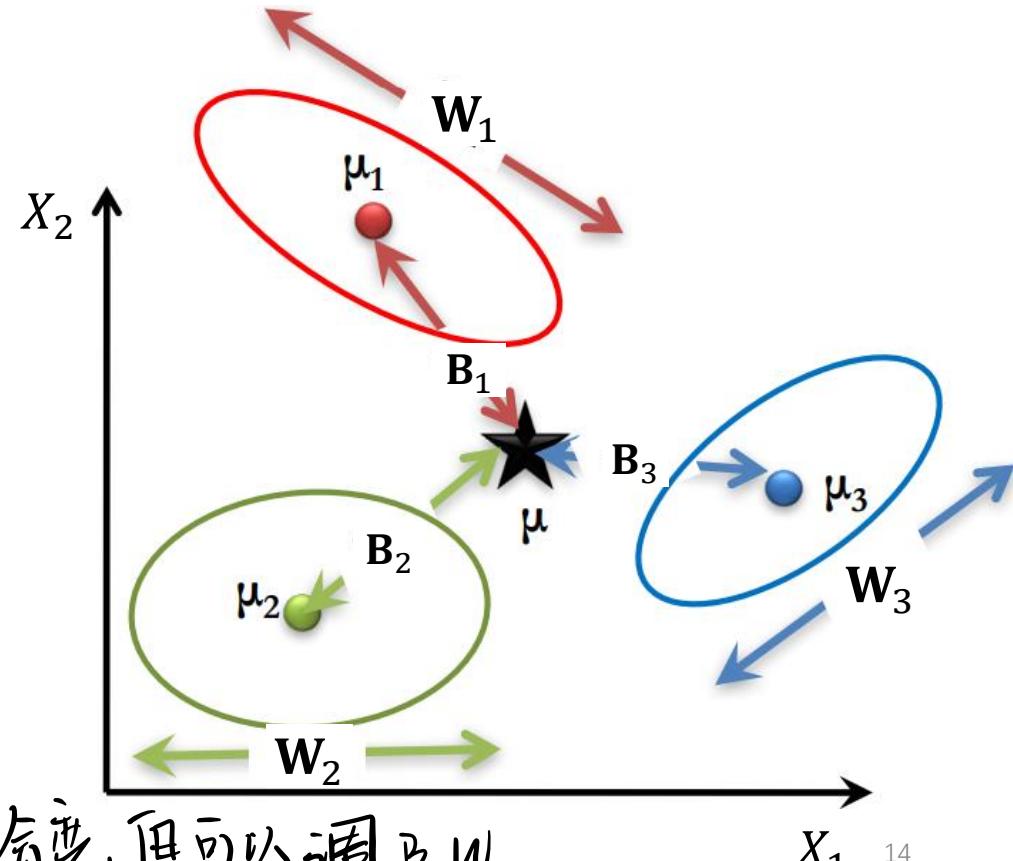
- Between class variance

$$\mathbf{B} = \sum_{k=1}^K N_k (\mu_k - \bar{\mu}) (\mu_k - \bar{\mu})^T$$

- Within class variance

$$\mathbf{W} = \sum_{k=1}^K \sum_{g_i=k} (x_i - \mu_k) (x_i - \mu_k)^T$$

给定数据非集后, $\mathbf{T} = \mathbf{B} + \mathbf{W}$ 不会变, 但可以调 \mathbf{B}, \mathbf{W}



Fisher's Formulation of Discriminant Analysis

- Maximize the Rayleigh quotient:

$$\max_a \frac{a^T \mathbf{B} a}{a^T \mathbf{W} a}$$

- Between class variance

$$\mathbf{B} = \sum_{k=1}^K N_k (\mu_k - \bar{\mu}) (\mu_k - \bar{\mu})^T$$

特征向量

- Within class variance

$$\mathbf{W} = \sum_{k=1}^K \sum_{g_i=k} (x_i - \bar{\mu}_k) (x_i - \bar{\mu}_k)^T$$

(假设W有逆)
 $\mathbf{W}^{-1} \mathbf{B} a = \lambda a$
 是对称的

- Equivalently, 约束优化
 \Rightarrow 拉格朗日
- $$\max_a a^T \mathbf{B} a$$
- $$\text{s.t. } a^T \mathbf{W} a = 1$$

$$L(a, \lambda) = a^T \mathbf{B} a - \lambda (a^T \mathbf{W} a)$$

- a is discriminant coordinates (canonical variates)

$$\frac{\partial L}{\partial a} = \mathbf{B} a - \lambda \mathbf{W} a = 0$$

- Generalized eigenvalue problem

$$\mathbf{B} a = \lambda \mathbf{W} a \quad \text{考虑 } \Delta x = \lambda x$$

which can be efficiently solved

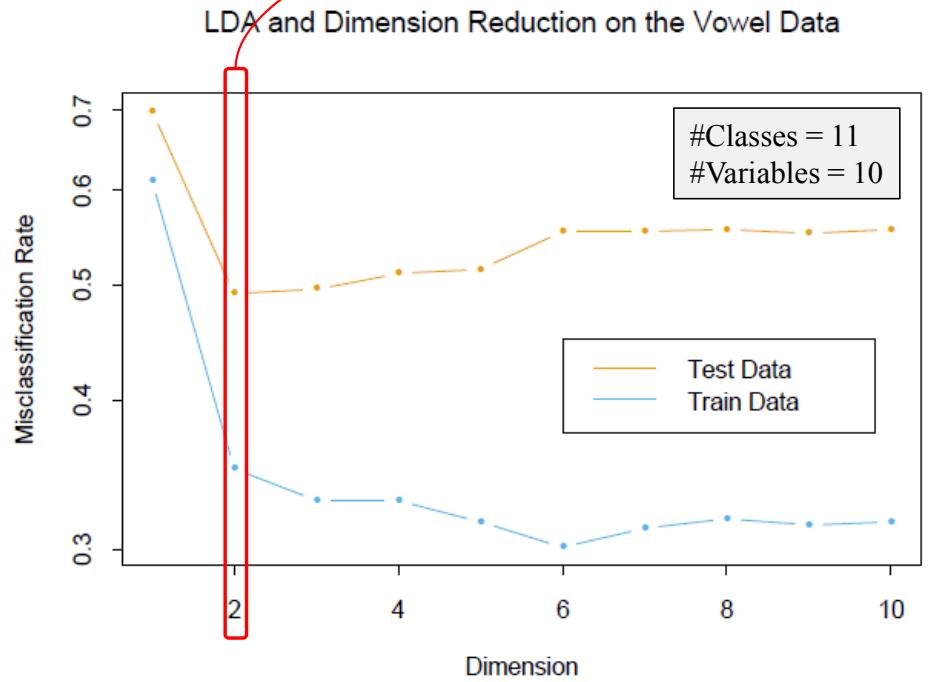
Ex. 4.1.

Hint: Lagrangian multipliers

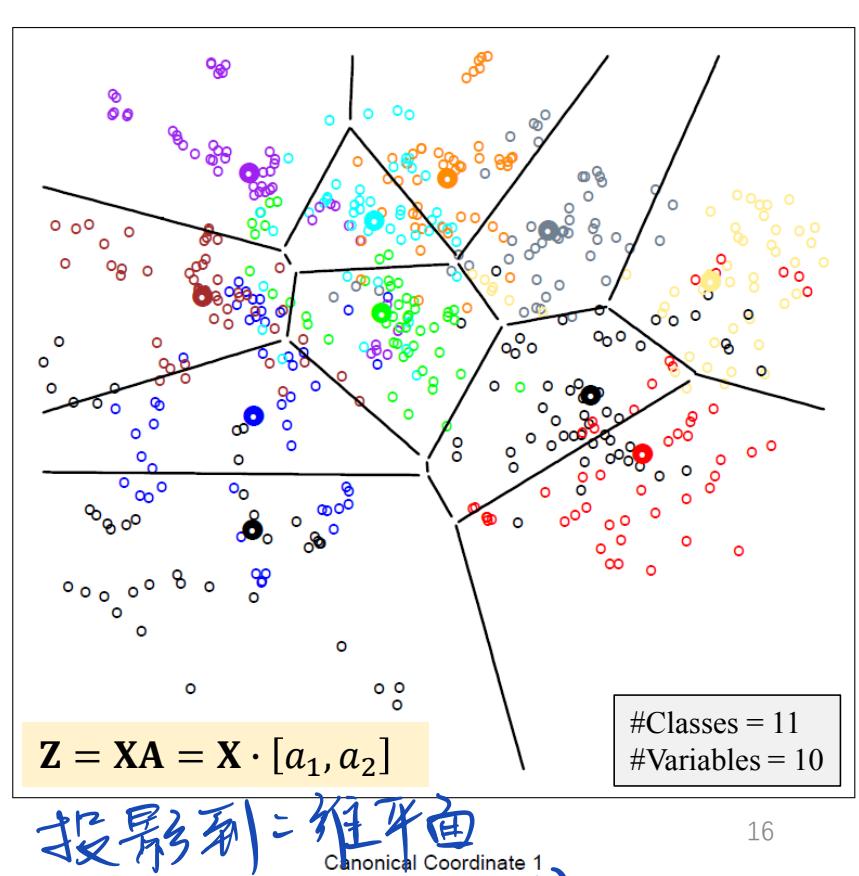
是用向量 a 来表示一个点， $a^T x_i$ 表示该点到 a 的距离。求 $z_i = [a_1 a_2 \dots a_k]^T f_i$

然后有 $3a$ ，计算投影 $z_i = f_i^T a$ ，考虑每个点的距离。

Fisher's Formulation of Discriminant Analysis



$$\max_a \frac{a^T \mathbf{B} a}{a^T \mathbf{W} a}$$



(有一些特征是无效的)

Linear Methods for Classification II

LDA

$$P(G_i=k|X=x) = \frac{P(X=x|G_i=k) P(G_i=k)}{P(X=x)}$$

→ 估计出后可以采样

$$P(X, Y)$$

如: LDA, NB

Native Bayes

- Generalization of LDA ~~似然式~~ Generative
- Logistic Regression 判别式方法: discriminative
- Summary

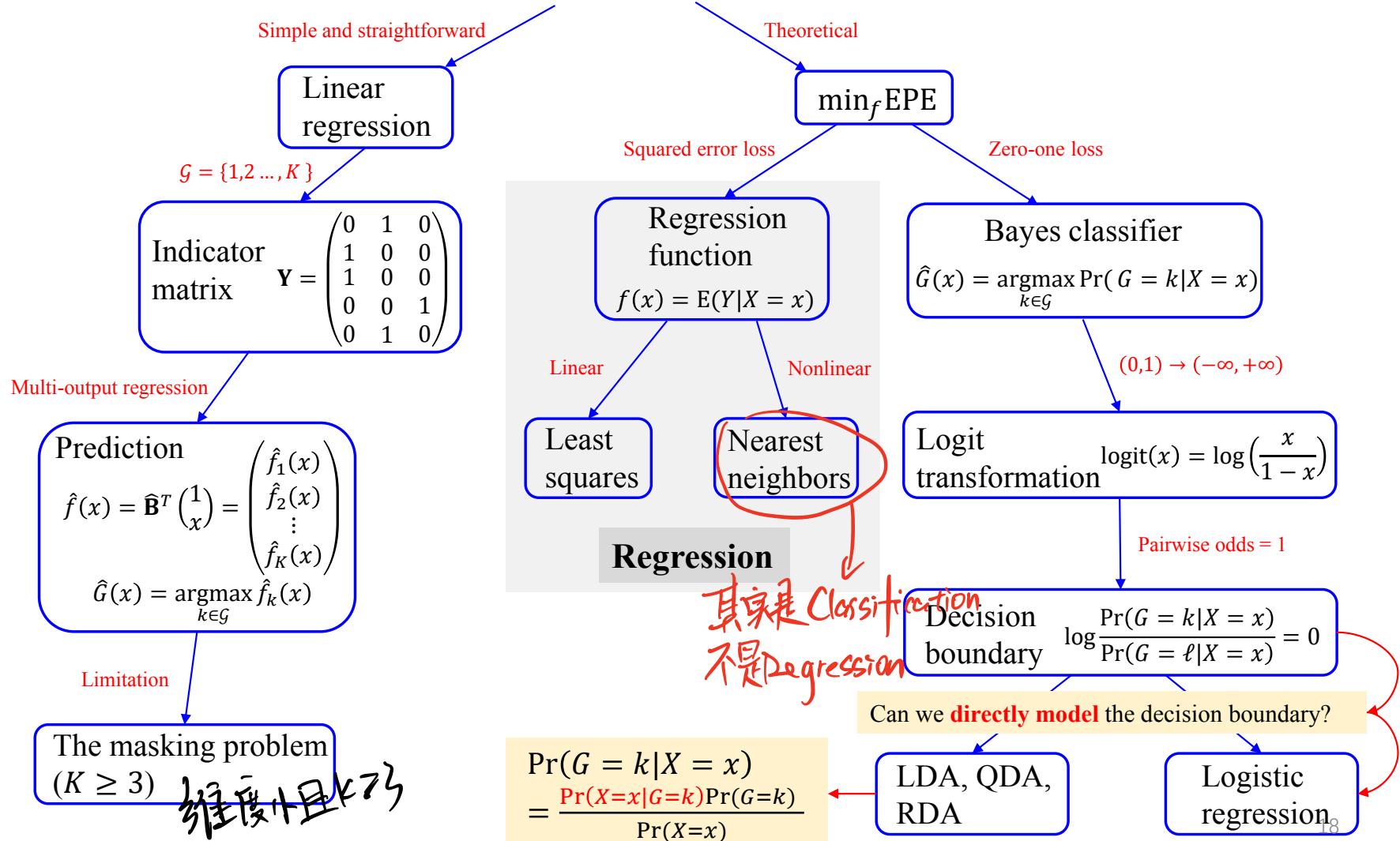
$$P(Y|X)$$

Ridge, Logistic, LS

岭回归

Least Square.

Classification



Linear Discriminant Analysis

- Recall our discussion on linear regression of an indicator matrix

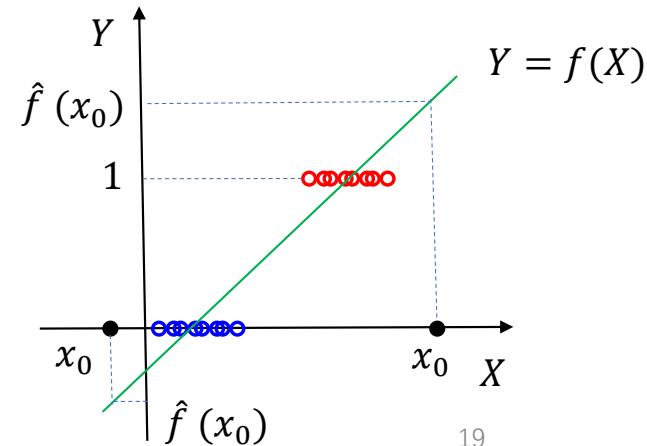
$$\hat{G}(x) = \operatorname{argmax}_{k \in \mathcal{G}} \hat{f}_k(x)$$

X

$$\hat{G}(x) = \operatorname{argmax}_{k \in \mathcal{G}} \Pr(G = k | X = x)$$

- It is inappropriate to represent a posterior directly by a linear function.
- Solution: make some **monotone transformation** of the posterior be linear in X

Linear decision boundary



Linear Discriminant Analysis

- Logit transform

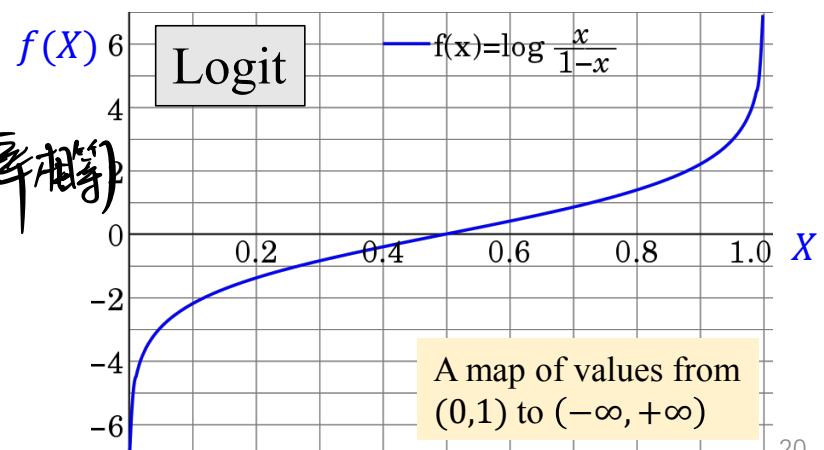
$$\text{logit}(\Pr(x)) = \log \left(\frac{\Pr(x)}{1 - \Pr(x)} \right)$$

It maps $\Pr(x) \in (0,1)$ to $\text{logit}(\Pr(x)) \in (-\infty, +\infty)$

- Decision boundary

- Odds equals to 1 →
发生 = 不发生
故在边界上 (概率相等)
- Or, logit equals to 0

Odds (发生比)



Linear Discriminant Analysis

Sigmoid 函数 $\sigma(x) = \frac{1}{1+e^x}$ 将 $(0, +\infty)$ 映射到 $(0, 1)$
 (也是 Logistic 函数)

- Example: binary (two class) classification
↑ 后验概率.

Logit: $\log \frac{\Pr(G=1|X=x)}{1-\Pr(G=1|X=x)} = \log \frac{\Pr(G=1|X=x)}{\Pr(G=2|X=x)} = \beta_0 + x^T \beta$

- The posterior probability

$$\Pr(G = 1|X = x) = \frac{\exp(\beta_0 + x^T \beta)}{1 + \exp(\beta_0 + x^T \beta)}, \quad \exp(x) = e^x$$

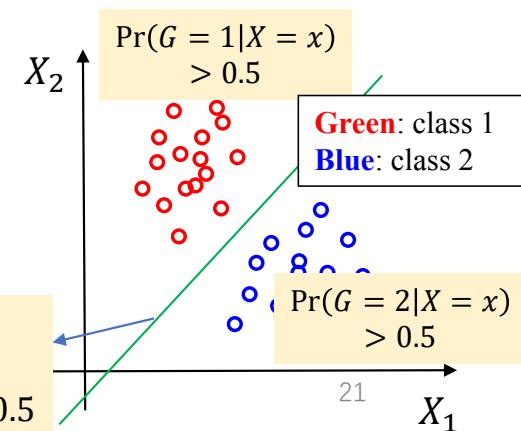
$$\Pr(G = 2|X = x) = \frac{1}{1 + \exp(\beta_0 + x^T \beta)}$$

- Decision boundary

$$\{x | \beta_0 + x^T \beta = 0\}$$

$$\sigma(x) = 1 - \sigma(x) \quad \sigma'(x) = \sigma(x)(1 - \sigma(x)) \text{ 但导数太小}$$

Decision boundary
 $\Pr(G = 1|X = x) = \Pr(G = 2|X = x) = 0.5$





Linear Logistic Regression

- Model the **posterior probabilities** of the K classes via linear function in x .

$$\log \frac{\Pr(G = 1|X = x)}{\Pr(G = K|X = x)} = \beta_{10} + x^T \beta_1$$

$$\log \frac{\Pr(G = 2|X = x)}{\Pr(G = K|X = x)} = \beta_{20} + x^T \beta_2$$

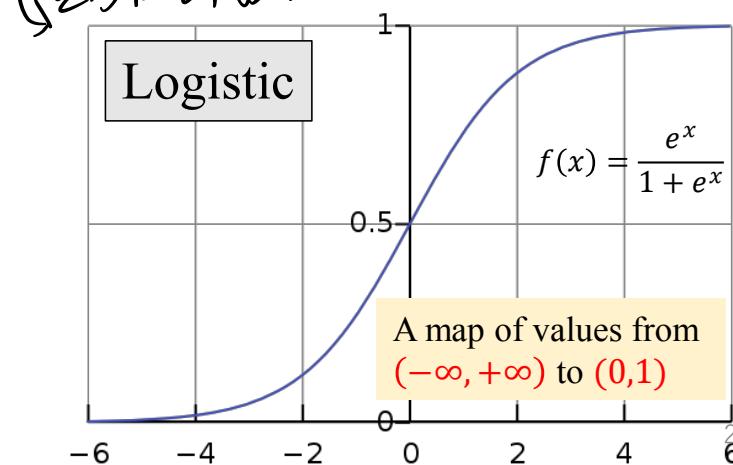
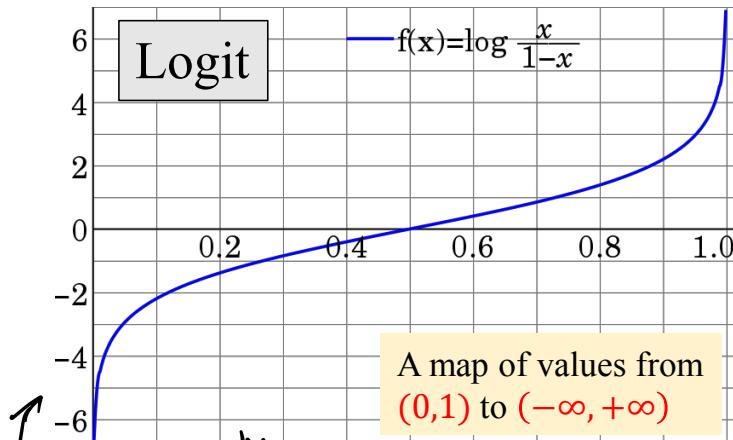
•

$$\log \frac{\Pr(G = K - 1 | X = x)}{\Pr(G = K | X = x)} = \beta_{(K-1)0} + x^T \beta_{K-1}$$

- $K - 1$ log-odds or **logit** function

$$\text{logitPr}(x) = \log \frac{\Pr(x)}{1 - \Pr(x)}$$

- The inverse of logit is **logistic** function



Linear Logistic Regression

- Model the **posterior probabilities** of the K classes via linear function in x .

$$\log \frac{\Pr(G = 1|X = x)}{\Pr(G = K|X = x)} = \beta_{10} + x^T \beta_1$$

$$\log \frac{\Pr(G = 2|X = x)}{\Pr(G = K|X = x)} = \beta_{20} + x^T \beta_2$$

:

$$\log \frac{\Pr(G = K - 1|X = x)}{\Pr(G = K|X = x)} = \beta_{(K-1)0} + x^T \beta_{K-1}$$

- $K - 1$ log-odds or **logit** function

$$\text{logitPr}(x) = \log \frac{\Pr(x)}{1 - \Pr(x)}$$

- The inverse of logit is **logistic** function

- A simple calculation yields

$$\Pr(G = k|X = x) = \frac{\exp(\beta_{k0} + x^T \beta_k)}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell0} + x^T \beta_\ell)}, \quad k = 1, \dots, K - 1$$

$$\Pr(G = K|X = x) = \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell0} + x^T \beta_\ell)}$$

- Parameter set

$$\theta = \{\beta_{10}, \beta_1, \dots, \beta_{(K-1)0}, \beta_{K-1}\}$$

- #parameters = $(p + 1) \times (K - 1)$



Linear Logistic Regression

$$\log \frac{\Pr(G = 1|X = x)}{\Pr(G = K|X = x)} = \beta_{10} + x^T \beta_1$$

⋮

$$\log \frac{\Pr(G = K - 1|X = x)}{\Pr(G = K|X = x)} = \beta_{(K-1)0} + x^T \beta_{K-1}$$

Pr($G = 1|X = x$) = Pr($G = K|X = x$) exp($\beta_{10} + x^T \beta_1$)

⋮

Pr($G = K - 1|X = x$) = Pr($G = K|X = x$) exp($\beta_{(K-1)0} + x^T \beta_{K-1}$)

summation

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

$\Pr(G = K|X = x) = \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + x^T \beta_{\ell})}$

$\Pr(G = k|X = x) = \frac{\exp(\beta_{k0} + x^T \beta_k)}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + x^T \beta_{\ell})}, k = 1, \dots, K - 1$

Linear Logistic Regression

- Estimating parameter set $\theta = \{\beta_{10}, \beta_1, \dots, \beta_{(K-1)0}, \beta_{K-1}\}$
 - Maximum likelihood estimation (MLE)
- Log-likelihood for N observations

$$\ell(\theta) = \log \Pr(\mathbf{g}|\mathbf{X}; \theta) = \sum_{i=1}^N \log \Pr(g_i|x_i; \theta)$$

- Two classes

$$\Pr(g=y|x; \theta) = p(x; \theta)^y (1 - p(x; \theta))^{1-y}$$

- Bernoulli distribution

- $\Pr(g=y|x; \theta) = p(x; \theta)^y (1 - p(x; \theta))^{1-y}$

Class	$g = 1$	$g = 2$
Code	$y = 1$	$y = 0$
Probability	$p(x; \theta)$	$1 - p(x; \theta)$

熵： $H(p) = -\int p(x) \log p(x) dx$ 反之熵： $H(p, q) = -\int p(x) \log q(x) dx$.

Linear Logistic Regression

- Two classes

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^N \{y_i \log p(x_i; \theta) + (1 - y_i) \log(1 - p(x_i; \theta))\} \\ &= \sum_{i=1}^N \left\{ y_i \left[x^T \beta - \log \left(1 + e^{x_i^T \beta} \right) \right] - (1 - y_i) \log \left(1 + e^{x_i^T \beta} \right) \right\} \\ &= \sum_{i=1}^N \left\{ y_i x_i^T \beta - \log \left(1 + e^{x_i^T \beta} \right) \right\} \end{aligned}$$

$x_i \leftarrow \begin{pmatrix} 1 \\ x_i \\ \beta_0 \\ \beta \end{pmatrix}$

没有闭式解 \Rightarrow 只能逼近

这里取 cross entropy. $p(x; \theta)$ 指观测的 y
 $q(x)$ 是计算得的

Please refer to:

https://en.wikipedia.org/wiki/Cross_entropy#Cross-entropy_loss_function_and_logistic_regression

$$\text{对数似然} : \textcircled{1} \ln \frac{\Pr(G_i=1|X=x)}{\Pr(G_i=2|X=x)} = 0 = \beta^T X$$

$$\textcircled{2} \Pr(G_i=1|X=x) = \frac{1}{1+e^{-\beta^T X}} = \Gamma(\beta_1^T X)$$

$$\Pr(G_i=2|X=x) = \Gamma(\beta_2^T X) = 1 - \Gamma(\beta_1^T X) = \Gamma(-\beta^T X)$$

$$\textcircled{3} \Pr(Y=y|X=x) = \Gamma^y(\beta^T X) (1 - \Gamma(\beta^T X))^{1-y}$$

$$\textcircled{4} \text{MLE} \max_{\beta} \ell(\beta) = \sum_{i=1}^n \ln \Pr(Y_i|X_i)$$

若 $y \in \{1, -1\}$ 时 可简化：\textcircled{3} \Pr(Y=y|X=x) = \Gamma(y\beta^T X)

$$\textcircled{4} \text{MLE} : \ell(\beta) = \sum_{i=1}^n \ln \frac{1}{1+e^{-y\beta^T X}} = -\sum_{i=1}^n \ln(1+e^{-y\beta^T X})$$

Linear Logistic Regression *

- The *first* derivative of $\ell(\theta)$

$$\begin{aligned}\frac{\partial \ell(\beta)}{\partial \beta} &= \sum_{i=1}^N \left(y_i x_i - \frac{\exp(x^T \beta)}{1 + \exp(x^T \beta)} \right) \\ &= \sum_{i=1}^N x_i (y_i - p(x_i))\end{aligned}$$

- The *second* derivative (Hessian)

$$\begin{aligned}\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} &= \sum_{i=1}^N -x_i \left(\frac{\partial p(x_i)}{\partial \beta^T} \right) \\ &= -\sum_{i=1}^N x_i x_i^T p(x_i) (1 - p(x_i))\end{aligned}$$

- In matrix form

$$\frac{\partial \ell(\beta)}{\partial \beta} = \mathbf{X}^T (\mathbf{y} - \mathbf{p})$$

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} = -\mathbf{X}^T \mathbf{W} \mathbf{X}$$

where $\mathbf{W} \in \mathbb{R}^{N \times N}$ is a diagonal matrix with the i -th diagonal element $p(x_i)(1 - p(x_i))$

The Newton-Raphson algorithm:

find the minimum or maximum iteratively by

$$x^{\text{new}} = x^{\text{old}} - \frac{f'(x^{\text{old}})}{f''(x^{\text{old}})}$$

- The Newton-Raphson step:

$$\begin{aligned}\beta^{\text{new}} &= \beta^{\text{old}} - \left(\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \ell(\beta)}{\partial \beta} \\ &= \beta^{\text{old}} + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \beta^{\text{old}} + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p})) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}\end{aligned}$$

- Given the response

$$\mathbf{z} = \mathbf{X} \beta^{\text{old}} + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p}),$$

- it is represented as a weighted least squares problem:

$$\beta^{\text{new}} \leftarrow \operatorname{argmin}_{\beta} (\mathbf{z} - \mathbf{X} \beta)^T \mathbf{W} (\mathbf{z} - \mathbf{X} \beta)$$

Linear Logistic Regression *

- Iteratively reweighted least squares (IRLS) algorithm

迭代加权最小二乘法。
随机或用 Least Square 简单计算一次。

1. Initialize β

2. Repeat

3. Form linearized responses

$$z_i = x_i^T \beta + \frac{y_i - p_i}{p_i(1 - p_i)} \quad \leftarrow \mathbf{z} = \mathbf{X}\beta^{\text{old}} + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p})$$

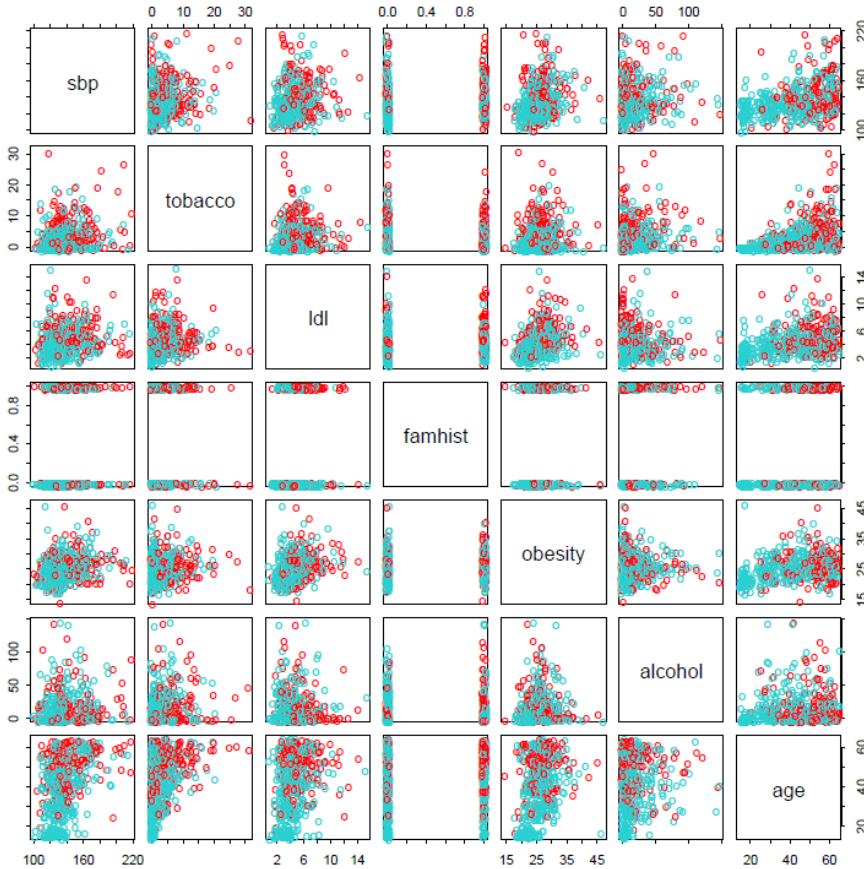
4. Form weights $w_i = p_i(1 - p_i)$

5. Update β by weighted least squares of z_i on x_i with w_i , $\forall i$

6. Until convergence

$$\beta^{\text{new}} \leftarrow \operatorname{argmin}_{\beta} (\mathbf{z} - \mathbf{X}\beta)^T \mathbf{W} (\mathbf{z} - \mathbf{X}\beta)$$

Linear Logistic Regression



Example: South African Heart Disease

- Red: 160 cases
- Green: 302 controls
- Z score measures the significance of a coefficient

	Coefficient	Std. Error	Z Score
sbp	0.006	0.006	1.023
tobacco	0.080	0.026	3.034
ldl	0.185	0.057	3.219
famhist	0.939	0.225	4.178
obesity	-0.035	0.029	-1.187
alcohol	0.001	0.004	0.136
age	0.043	0.010	4.184

收缩压

肥胖
饮酒

The data is fitted by logistic regression

Linear Logistic Regression

- L_1 regularized logistic regression

$$\max_{\beta_0, \beta} \left\{ \sum_{i=1}^N \left[y_i (\beta_0 + \beta^T x_i) - \log(1 + e^{\beta_0 + \beta^T x_i}) \right] - \lambda \sum_{j=1}^p |\beta_j| \right\}$$

超参数.

$\hookrightarrow \text{正则化}$

- Standardize the inputs, and penalize without β_0
- Solved by the Newton algorithm
 - Replace the weighted least squares by the weighted lasso.
- L_2 regularized logistic regression? Algorithm?
 \hookrightarrow 加权岭回归

Connection between LDA and Logistic Regression

- The linear logistic model only specifies the **conditional distribution**, while the LDA model specifies the **joint distribution**
- If the additional assumption LDA 需要更多假设 made by LDA is appropriate, LDA tends to estimate the parameters more efficiently.
- Another advantage of LDA is that samples without class labels LDA 可以做无监督学习 can be used under the model of LDA. On the other hand, LDA is not robust to gross outliers. Because logistic regression relies on fewer assumptions, it seems to be more robust to the non-Gaussian type of data.
- In practice, logistic regression and LDA often give **similar results**.

Linear Methods for Classification II

- Generalization of LDA
 - Regularized Discriminant Analysis
 - Fisher's Formulation of Discriminant Analysis
- Logistic Regression
- Summary

