

Discussion 03

2022.03.10

杨子健 yangzi@shanghaitech.edu.cn

p : number of features

| | 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 \mathbf{X}

N: number of instances

| G | | Y_1 | Y_2 | \dots | Y_{10} |
|---------|------------------------|---------|---------|---------|----------|
| 0 | One-hot coding → | 1 | 0 | \dots | 0 |
| 1 | | 0 | 1 | \dots | 0 |
| \dots | | \dots | \dots | \dots | \dots |
| 9 | | 0 | 0 | \dots | 1 |

Coded output matrix \mathbf{Y}

N: number of instances

Q $\min_{\mathbf{B}} \|\mathbf{Y} - \mathbf{XB}\|_F^2 \rightarrow \hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

$$\|A\|_F^2 = \text{tr}(A^T A)$$

$$\begin{aligned} \|\mathbf{Y} - \mathbf{XB}\|_F^2 &= \text{tr}[(\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB})] \\ &= \text{tr}[\mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{XB} - \mathbf{B}^T \mathbf{X}^T \mathbf{Y} + \mathbf{B}^T \mathbf{X}^T \mathbf{XB}] \end{aligned}$$

$$\text{tr}(A) = \text{tr}(A^T)$$

$$= \text{tr}(\mathbf{Y}^T \mathbf{Y}) + \text{tr}(\mathbf{B}^T \mathbf{X}^T \mathbf{XB}) - 2\text{tr}(\mathbf{Y}^T \mathbf{XB})$$

$$\nabla_{\mathbf{B}} \|\mathbf{Y} - \mathbf{XB}\|_F^2 = \nabla_{\mathbf{B}} \text{tr}(\mathbf{B}^T \mathbf{X}^T \mathbf{XB}) - \nabla_{\mathbf{B}} 2\text{tr}(\mathbf{Y}^T \mathbf{XB})$$

$$\nabla_{\mathbf{X}} \text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{X}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{X} \quad \nabla_{\mathbf{X}} \text{tr}(\mathbf{A} \mathbf{X} \mathbf{C}) = \mathbf{A}^T \mathbf{C}^T$$

$$\nabla_{\mathbf{B}} \|\mathbf{Y} - \mathbf{XB}\|_F^2 = (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \mathbf{B} - 2 \mathbf{X}^T \mathbf{Y} \stackrel{!}{=} 0$$

$$\Rightarrow \hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

Binary classification

- Linear regression

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

- Decision boundary

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

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

Multi-class classification

- Linear regressions for K classes

$$f_k(x) = \beta_{k0} + x^T \beta_k, \quad k = 1, \dots, K$$

- Decision boundary between classes k and ℓ :

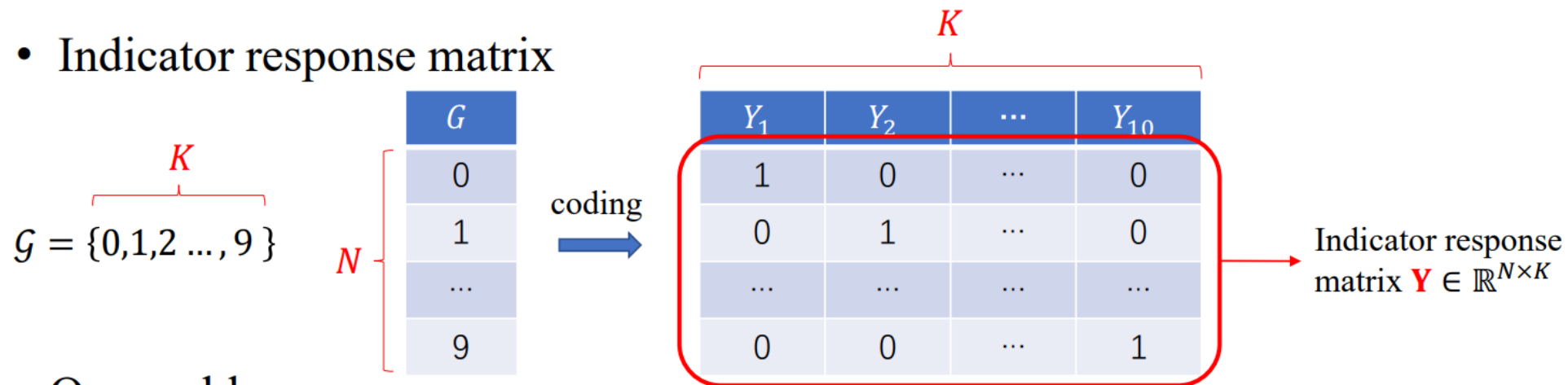
$$\{x : \hat{f}_k(x) = \hat{f}_\ell(x)\}$$

- 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 Regression of an Indicator Matrix

- Indicator response matrix



- Our problem:

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

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

- The fitted values on \mathbf{X} :

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

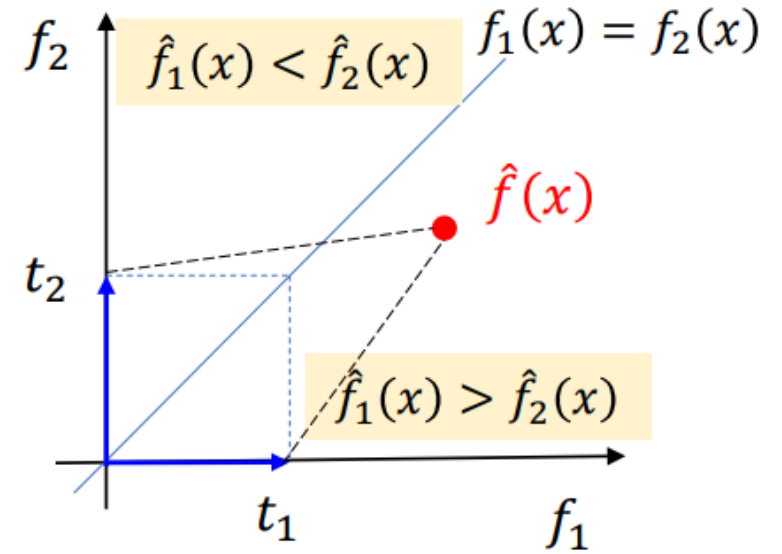
A new observation x is classified by

- Compute the fitted output

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

- Classify x according to

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



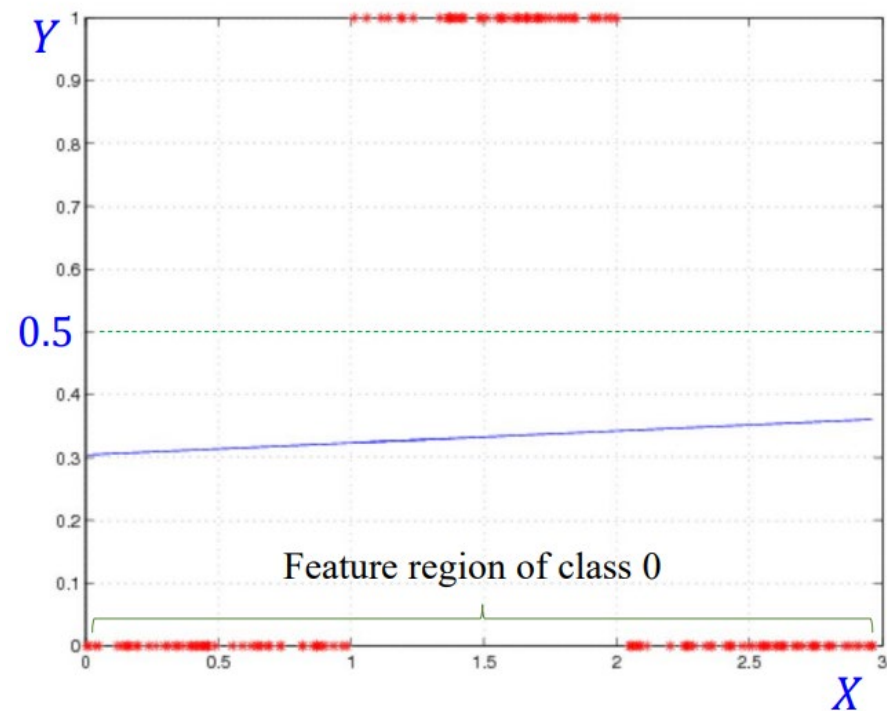
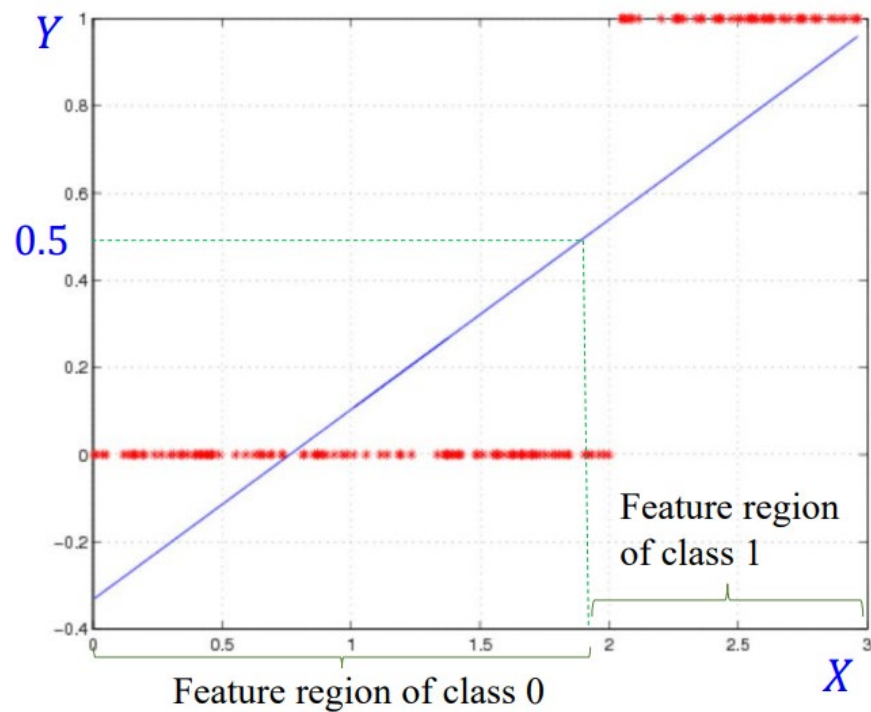
$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\mathbf{B}}$$

$$\underset{\mathcal{N}}{\overset{y_i^T}{\cancel{\boxed{\mathbf{Y}}}}}^k = \underset{\mathcal{N}}{\overset{x_i^T}{\cancel{\boxed{\mathbf{X}}}}}^{p+1} \cdot \overset{p+1}{\boxed{\hat{\mathbf{B}}}}^k$$

$$y_i^T = x_i^T \cdot \hat{\mathbf{B}}$$

$$y_i = \hat{\mathbf{B}}^T \cdot x_i = f(x_i)$$

The Phenomenon of Masking



Linear Discriminant Analysis

- **Example:** binary (two class) classification

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

Q

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

- Decision boundary

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

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

$\nearrow p_1$
 $\searrow p_2$
记为 t

$$\log \frac{p_1}{p_2} = t \Rightarrow \frac{p_1}{p_2} = e^t$$

$$p_1 + p_2 = (e^t + 1) p_2 = 1$$

$$\Rightarrow p_2 = \frac{1}{e^t + 1}, \quad p_1 = \frac{e^t}{e^t + 1}$$

- Posterior

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

□ Density of X in class $G = k$:

$$f_k(x) = \Pr(X = x|G = k)$$

□ Class prior:

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

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

- Assumptions in LDA

1. Model each class density as **multivariate Gaussian**

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

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

- Compare two classes k and ℓ

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

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

Quadratic term
vanished due to
the common
covariance

Quadratic Discriminant Analysis

- **Assumption:** Each class has a specific covariance Σ_k

- **Difference with LDA**

- Σ_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

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

Σ

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

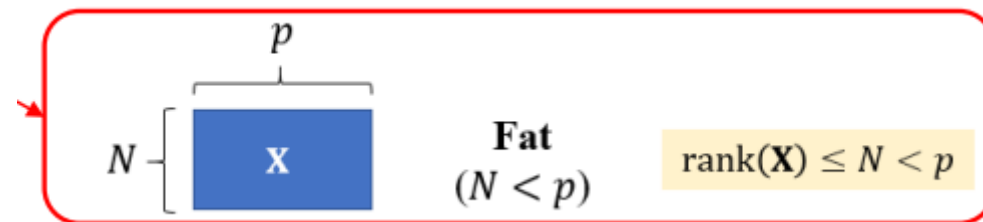
Regularized Discriminant Analysis

High dimensional problems ($p \gg N$)

- Regularization is necessary
 - No enough data to estimate feature dependencies
 - E.g., independent assumption on features
 - Diagonal within-class covariance matrix
- #paras: $K \times p \times p \rightarrow K \times p$

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 $\hat{\Sigma}$



Regularized Discriminant Analysis

Regularized Discriminant Analysis on the Vowel Data

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

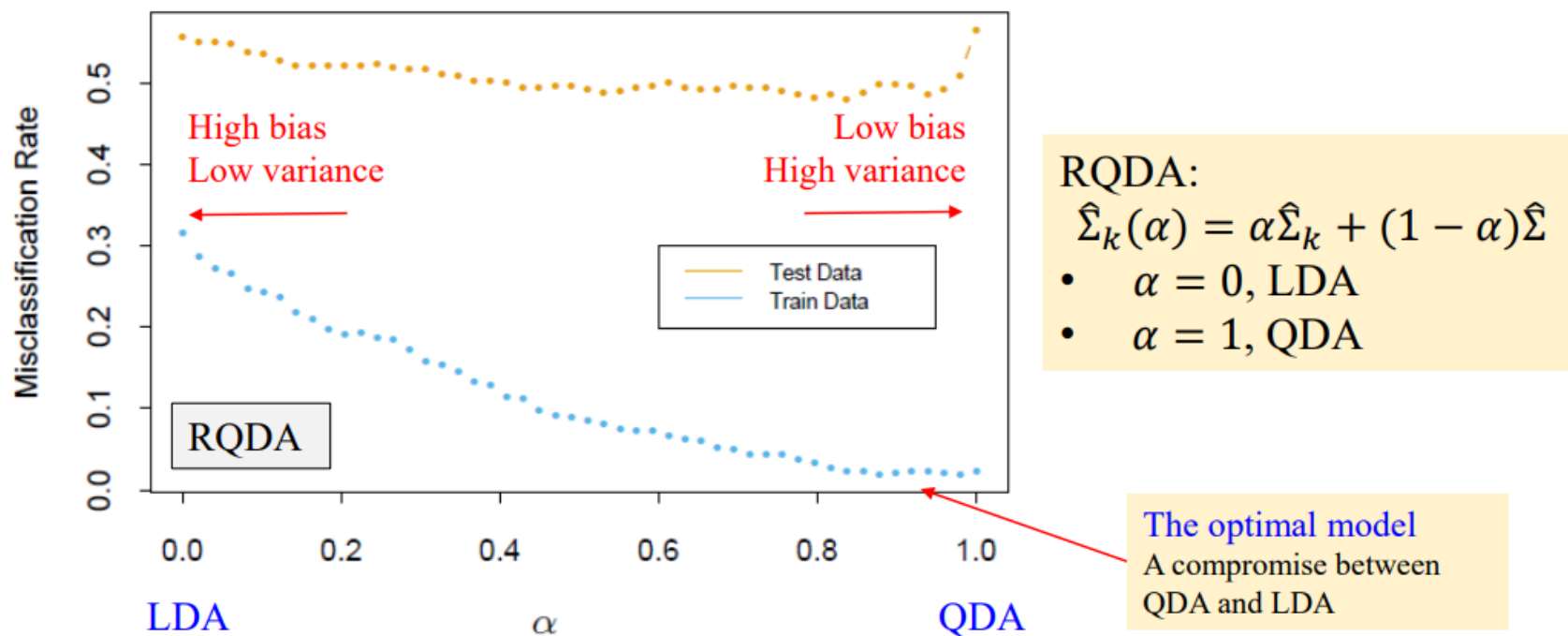
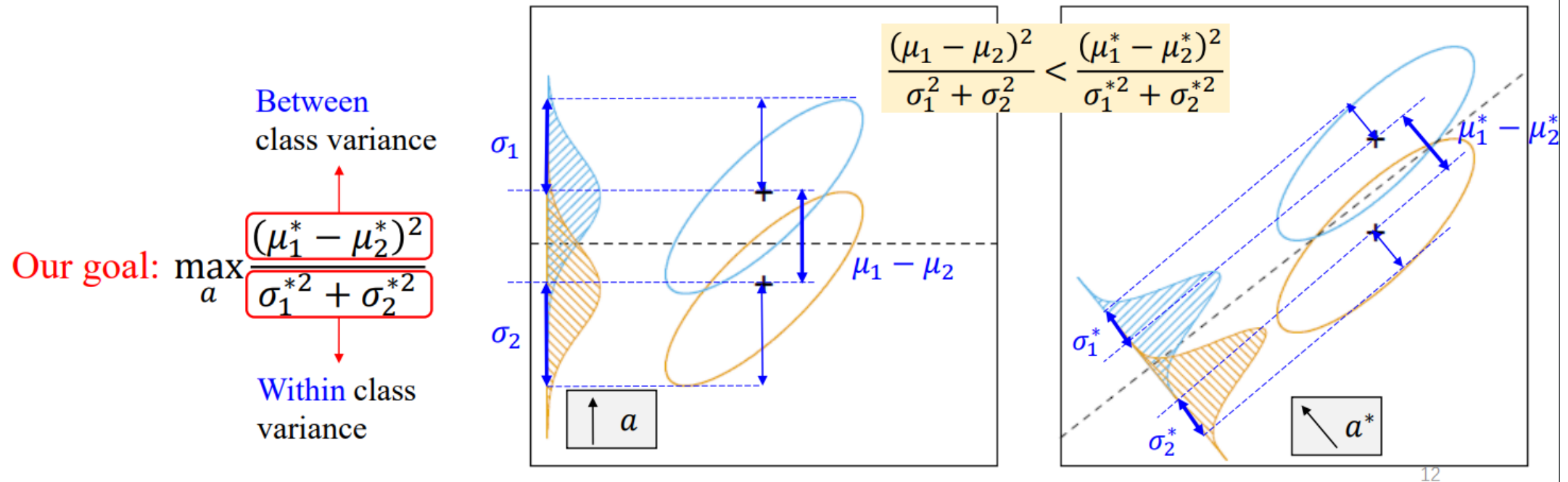


FIGURE 4.7. Test and training errors for the vowel data, using regularized discriminant analysis with a series of values of $\alpha \in [0, 1]$. The optimum for the test data occurs around $\alpha = 0.9$, close to quadratic discriminant analysis.

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.



THANKS