

---

# **CS722/822: Machine Learning**

Instructor: Jiangwen Sun  
Computer Science Department

---

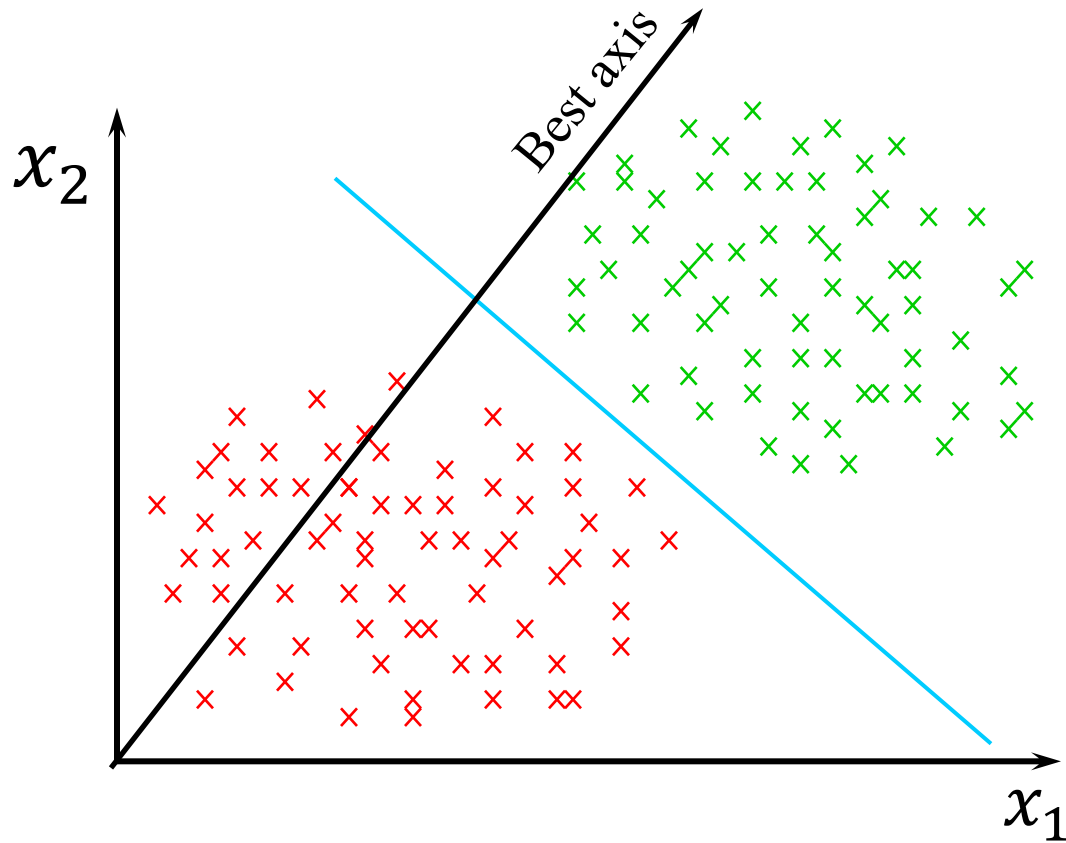
# **Linear Discriminant Analysis (Fisher's Linear Discriminant)**

# Classification

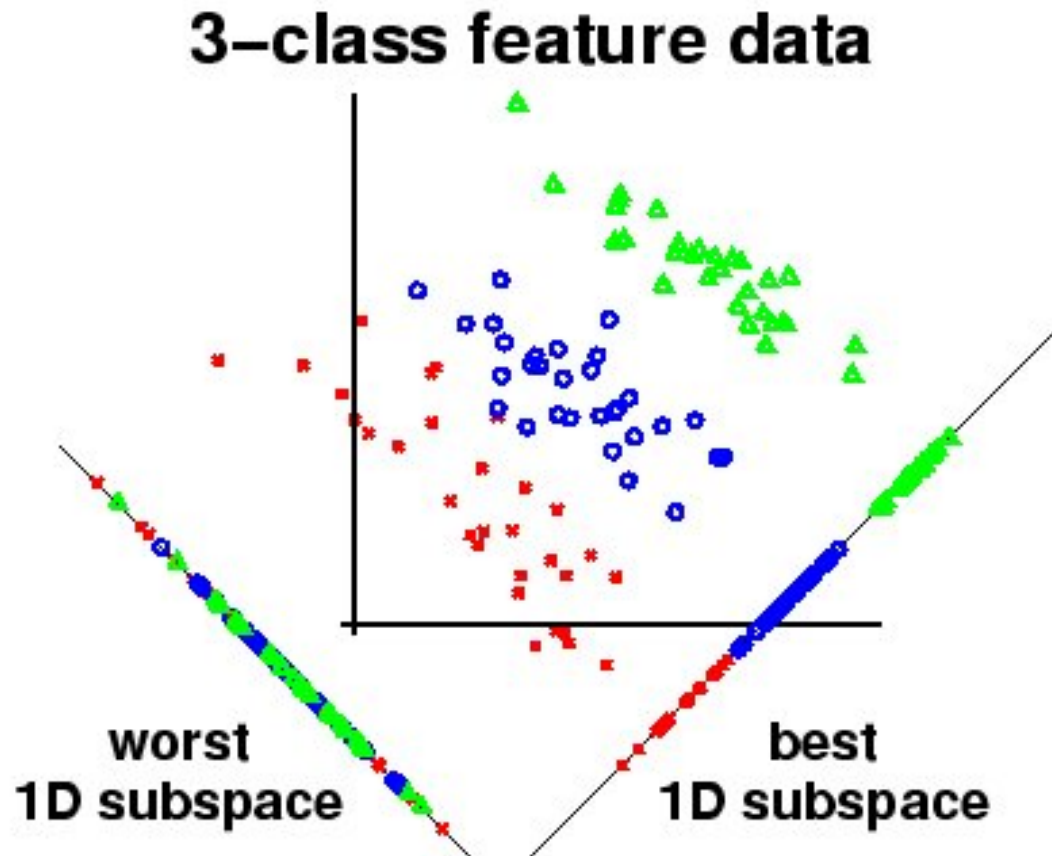
---

- Training data is given
  - Each object is associated with a class label  $y \in \{1, 2, \dots, K\}$  and a feature vector of  $d$  measurements:  $\mathbf{x} = (x_1, \dots, x_d)$ .
- Build a predictive model:  $f(\mathbf{x}) \rightarrow y$  from the training data.
- Unseen objects are to be classified as belonging to one of a number of predefined classes  $\{1, 2, \dots, K\}$ .
- Linear Discriminant Analysis / Fisher's linear discriminant

# Two classes



# Three classes



# Classifiers

---

- Given a training set

$$L = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

- Classifier  $f$  built from  $L$ :

$$f: x \rightarrow y \in \{1, 2, \dots, K\}$$

- **Bayes classifier** base on conditional densities  $p(k|x)$ ,

$$f(x) = \operatorname{argmax}_k p(k|x)$$

This is a *maximum a posterior*, and  $p(k|x)$  is a posterior density

# The Rules of Probability

- **Sum Rule**

$$p(X) = \sum_Y p(X, Y)$$

- **Product Rule**

$$\begin{aligned} p(X, Y) &= p(Y|X)p(X) \\ &= p(X|Y)p(Y) \end{aligned}$$

- **Bayes' Rule**

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

$p(X) = \sum_Y p(X|Y)p(Y)$  is irrelevant to  $Y = k$

$$p(Y = k|X = x) \propto p(X = x|Y = k)p(Y = k)$$

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

# Maximum a posterior

---

- $p(k|x) = p(x|k)p(k)/p(x)$

- Find a class label  $k$  so that

$$\operatorname{argmax}_k p(k|x) = \operatorname{argmax}_k p(x|k)p(k)$$

- Naïve Bayes assumes **independence** among all features given the class:

$$p(x|k) = p(x_1|k)p(x_2|k) \cdots p(x_d|k)$$

**Very strong assumption**



# Multivariate normal dist for each class

Assume multivariate Gaussian (normal) class densities  
 $(\mathbf{x}|k) \sim N(\mu_k, \Sigma_k)$ :

$$p(\mathbf{x}|k) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_k|}} \exp\left(-\frac{1}{2} \left( (\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k) \right)\right)$$


**Maximizing posterior** is equivalent to **maximizing  $p(\mathbf{x}|k)p(k)$** , and  
equivalent to **maximizing the logarithm of  $p(\mathbf{x}|k)p(k)$**


$$\begin{aligned} \log(p(\mathbf{x}|k)p(k)) &= -\frac{1}{2} \left( (\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k) \right) \\ &\quad - \log\left(\sqrt{(2\pi)^d |\Sigma_k|}\right) + \log(p(k)) \end{aligned}$$


$$f(\mathbf{x}) = \operatorname{argmin}_k \{ (\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k) + \log|\Sigma_k| - 2 \log(p(k)) \}$$


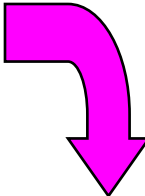
# Two-class case

Two class labels:  $k_1$  and  $k_2$

If  $p(x|k_1)p(k_1) > p(x|k_2)p(k_2)$    $f(x) = k_1$

otherwise   $f(x) = k_2$

Equivalently,  $\frac{p(x|k_1)p(k_1)}{p(x|k_2)p(k_2)} > 1$    $\frac{p(x|k_1)}{p(x|k_2)} > \frac{p(k_2)}{p(k_1)}$

  $\log \frac{p(x|k_1)}{p(x|k_2)} > \log \frac{p(k_2)}{p(k_1)}$  

$$(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \log |\Sigma_1| - (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) + \log |\Sigma_2| < T$$

$T$  is a constant

# Guassain discriminant rule

- For multivariate Gaussian (normal) class densities, i.e.,

$$\{x|(y = k)\} \sim N(\mu_k, \Sigma_k),$$

the classification rule (predictive function or model) is

$$f(x) = \underset{k}{\operatorname{argmin}} \{ (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log |\Sigma_k| - 2 \log(p(k)) \}$$

- In two-class, this is a **quadratic rule** (Quadratic discriminant analysis, or QDA)

$$(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \log |\Sigma_1| - (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) + \log |\Sigma_2| < T$$

- In practice, population mean vectors  $\mu_k$  and covariance matrices  $\Sigma_k$  are estimated by corresponding sample quantities

# Sample mean and variance

---

- Class mean

$$\mu_k = \frac{1}{|C_k|} \sum_{x \in C_k} x$$

- Class covariance

$$\Sigma_k = \frac{1}{|C_k|} \sum_{x \in C_k} (x - \mu_k)(x - \mu_k)^T$$

# Example

$$X_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \quad X_2 = \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix}, \quad X_3 = \begin{pmatrix} 0 \\ 2 \\ 1 \end{pmatrix}.$$

$$\mu = \frac{1}{3}(X_1 + X_2 + X_3) = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\Sigma = \frac{1}{3}((X_1 - \mu)(X_1 - \mu)^T + (X_2 - \mu)(X_2 - \mu)^T + (X_3 - \mu)(X_3 - \mu)^T)$$

$$= \frac{1}{3} \left( \begin{pmatrix} 0 \\ -1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 & -1 & 0 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \end{pmatrix} + \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} \begin{pmatrix} -1 & 1 & 0 \end{pmatrix} \right)$$

$$= \frac{1}{3} \left( \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 1 & -1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \right) = \frac{1}{3} \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

# Two-class case

- If the two classes have the same covariance matrix,  $\Sigma_k = \Sigma$ , the discriminant rule is **linear** (Linear discriminant analysis, or LDA; FLDA for  $K = 2$ ):
- Quadratic rule

$$(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \log |\Sigma_1| - (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) + \log |\Sigma_2| < T$$

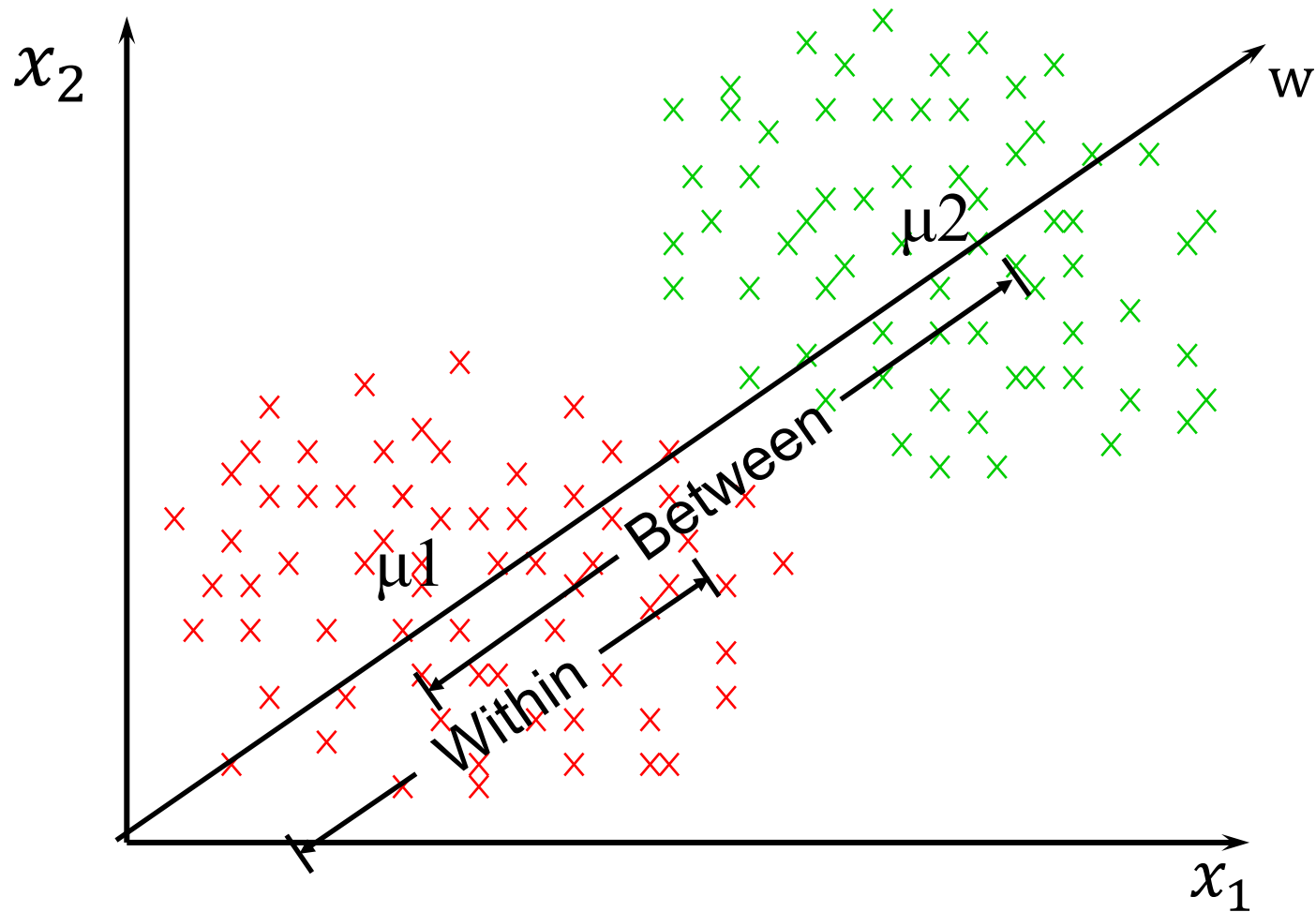
$$\text{Become } x^T \Sigma^{-1} (\mu_2 - \mu_1) < c$$

$$\boxed{x^T w < c}, \text{ where } w = \Sigma^{-1} (\mu_2 - \mu_1)$$

In practice, **Linear rule**

$$\Sigma = \frac{1}{n} (n_1 \Sigma_1 + n_2 \Sigma_2)$$

# Illustration



## Good separation

# Two-class case

- Maximize the signal-to-noise ratio

$$\max_w \frac{w^T \Sigma_{between} w}{w^T \Sigma_{within} w} \quad \begin{array}{l} \longrightarrow \text{Between-class separation} \\ \longrightarrow \text{Within-class separation} \end{array}$$

$s. t. \|w\| = 1$

where  $\Sigma_{between} = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T$

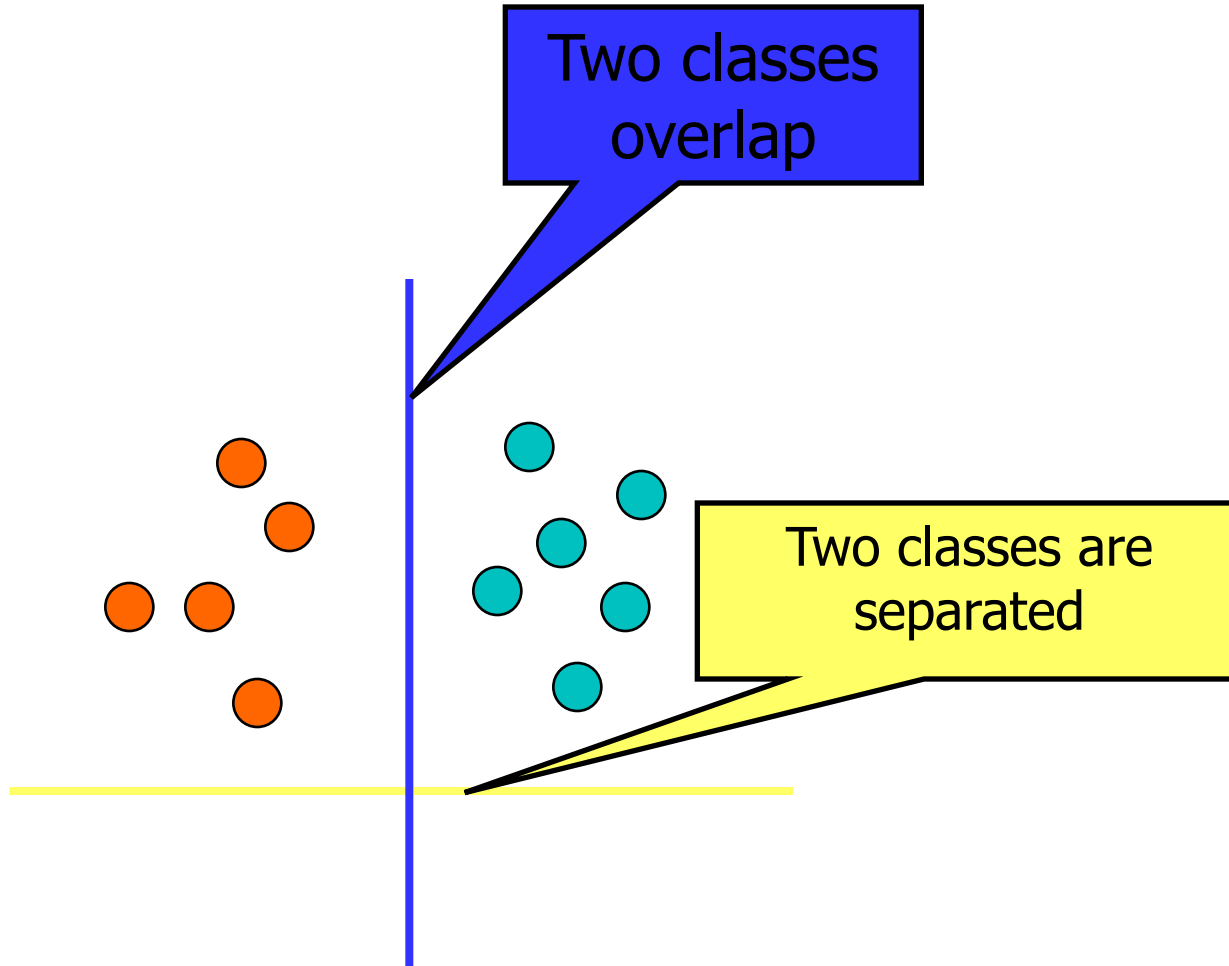
$$\Sigma_{within} = \frac{1}{n} (n_1 \Sigma_1 + n_2 \Sigma_2)$$

Solution is  $\Sigma_{within}^{-1}(\mu_2 - \mu_1)$

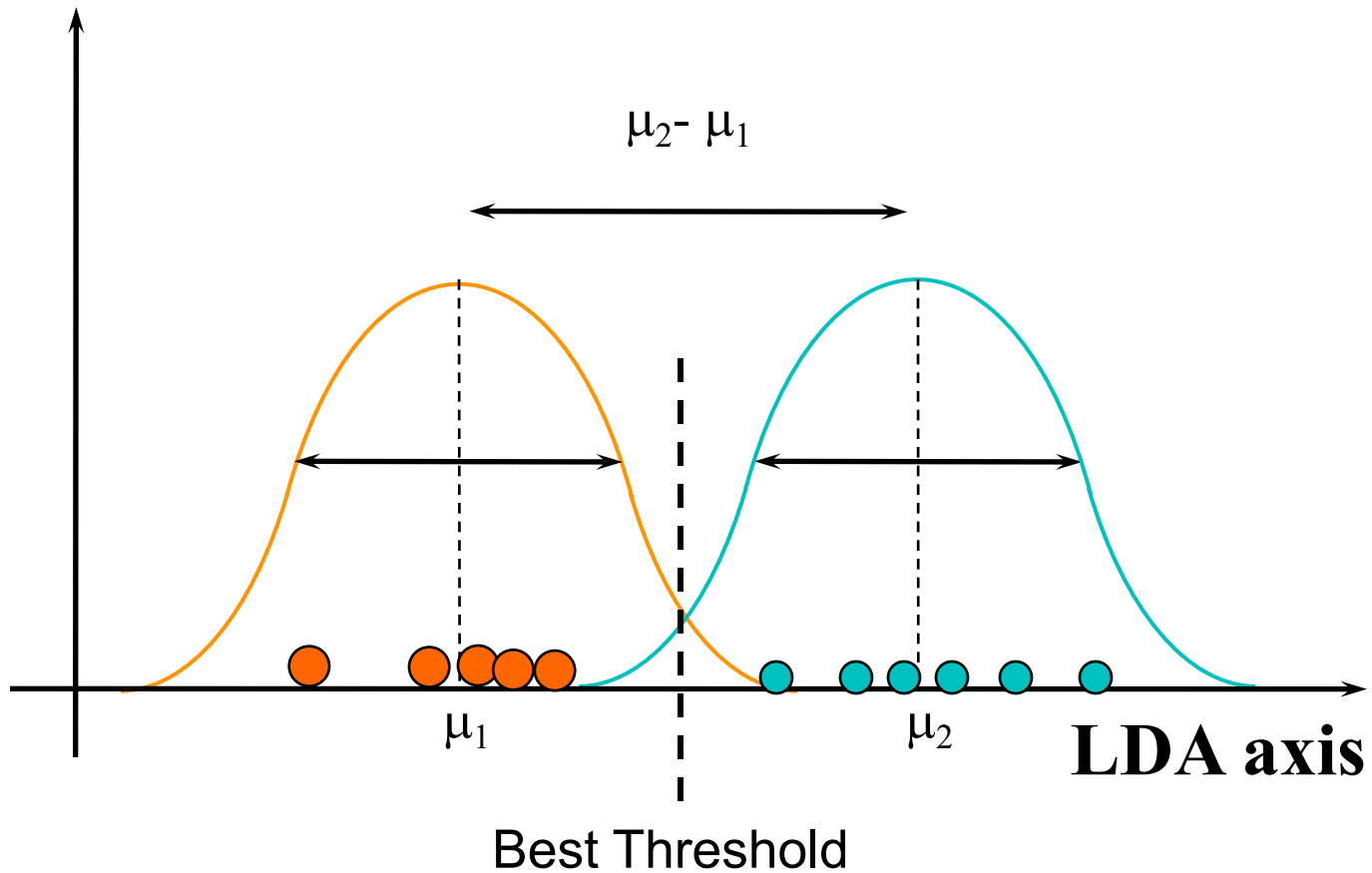


# Two-class case (illustration)

- LDA gives the yellow direction

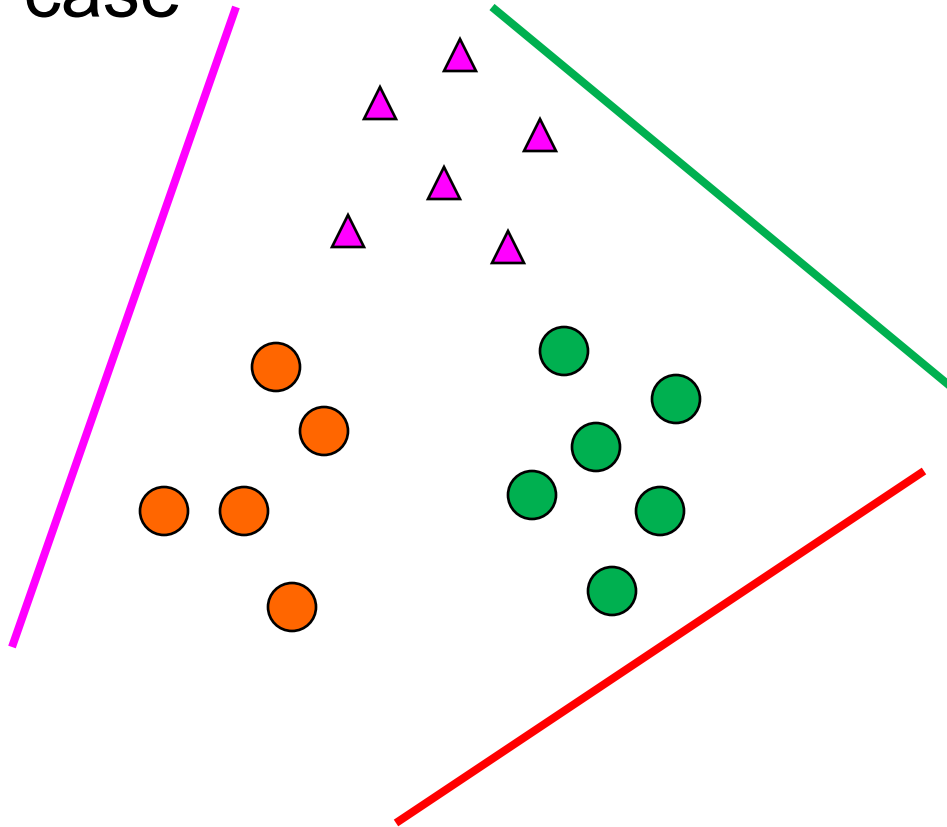


# Two-class case (illustration)



# Multi-class case

- Two approaches
  - Apply two-class LDA to each “one-versus-rest” case



# Multi-class case

---

**Second approach:** find multiple directions that form a low dimensional space

Transformation matrix  $W$  that projects the data to be most separable is the one that maximizes

$$\max_W \text{trace} \left( \frac{W^T S_b W}{W^T S_w W} \right)$$
$$s. t. W^T W = I$$

Between-class matrix:  $S_b = \frac{1}{n} \sum_{i=1}^K n_i (\mu_i - \mu)(\mu_i - \mu)^T$

Within-class matrix:  $S_w = \frac{1}{n} \sum_{i=1}^K \sum_{x \in C_i} (x - \mu_i)(x - \mu_i)^T$

# Intuition

---

The goal is to simultaneously maximize the between-class separation and minimize the within-class separation

The solution to:

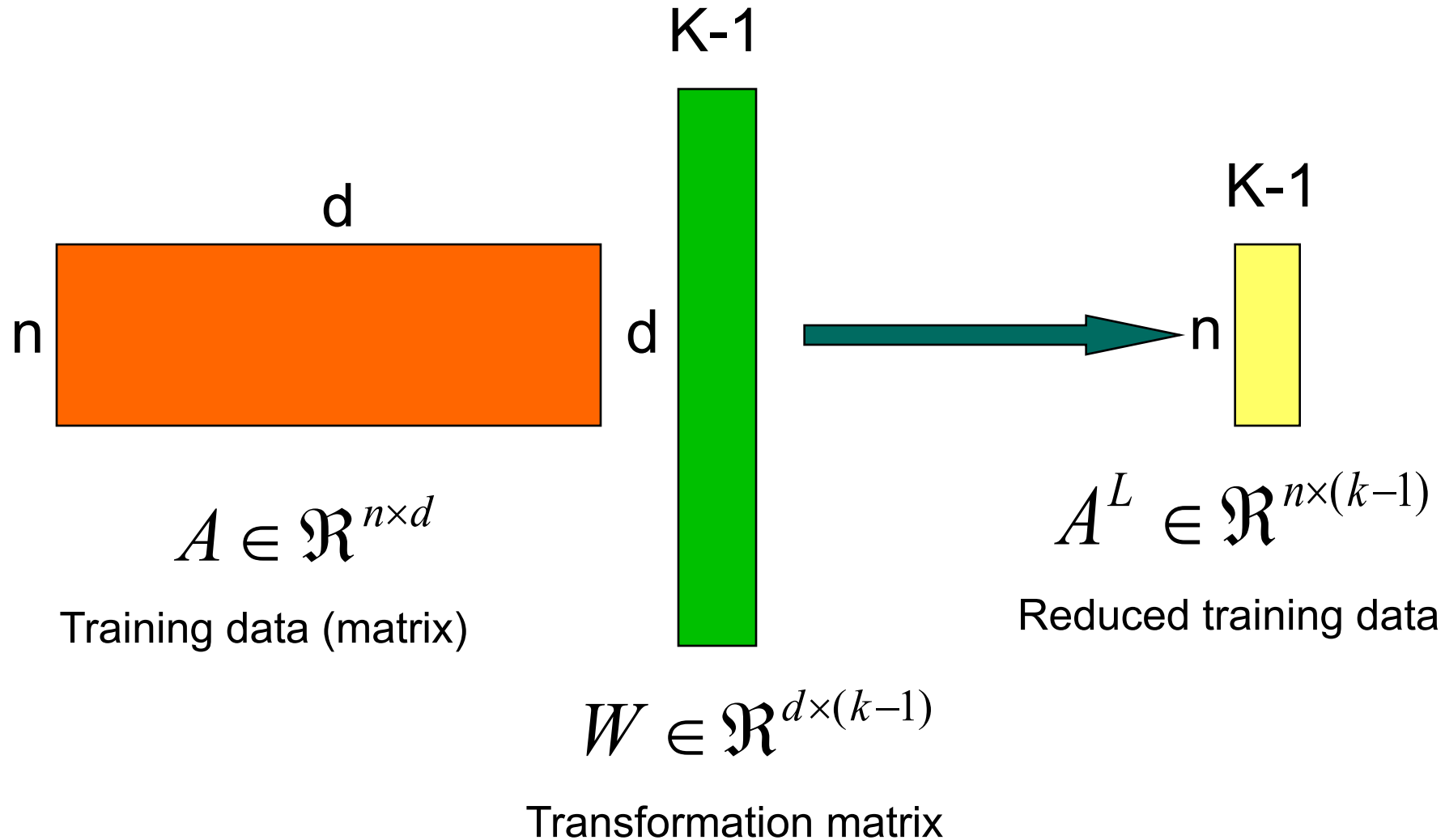
$$\begin{aligned} \max_W \text{trace} \left( \frac{W^T S_b W}{W^T S_w W} \right) \\ \text{s.t. } W^T W = I \end{aligned}$$

is the generalized eigenvalue problem:

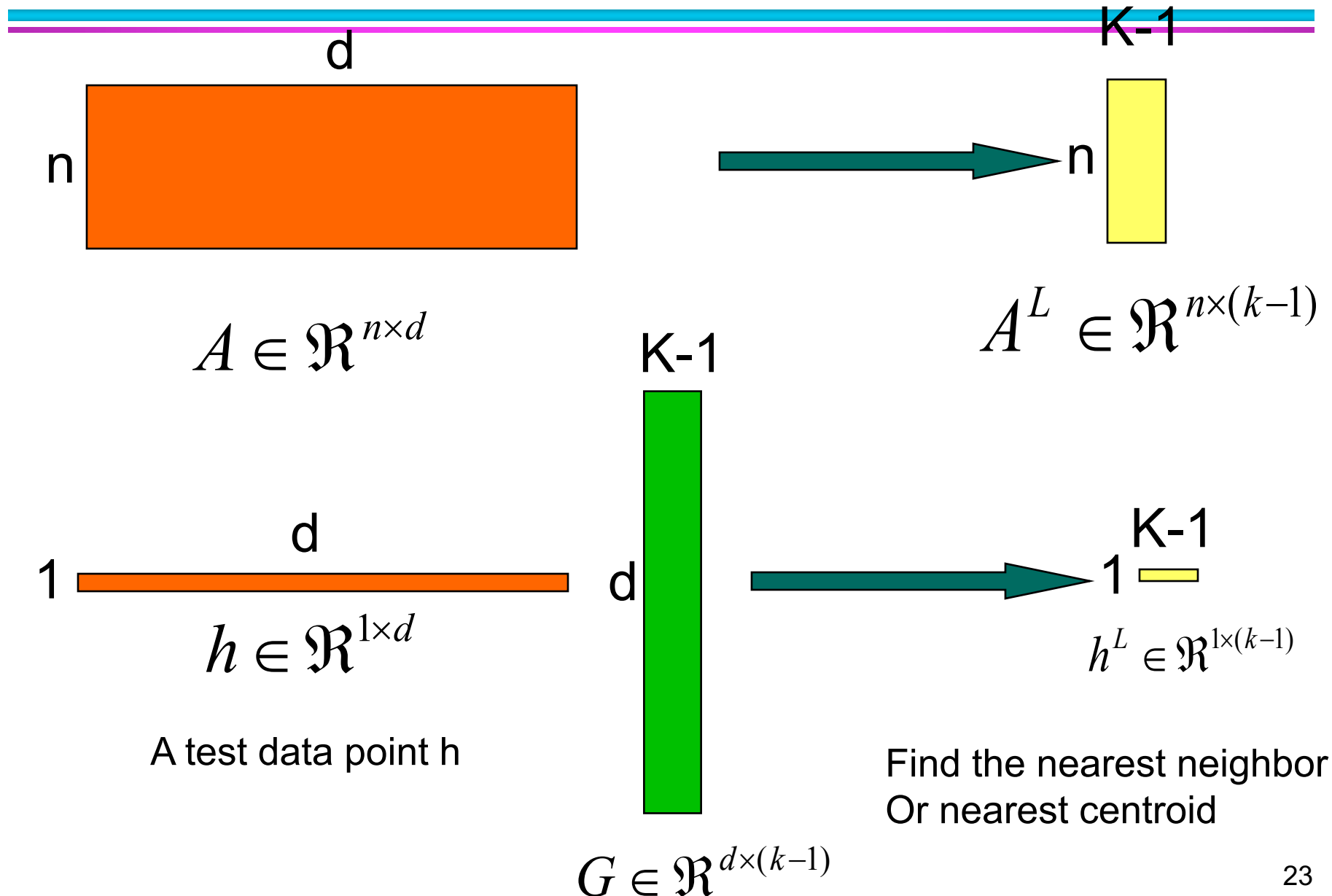
$$S_b w = \lambda S_w w$$

where we want to find the eigenvectors associated with the first  $k - 1$  largest eigenvalues.

# Graphic view of the transformation (projection)



# Graphical view of classification



# Summary

---

- Dimension reduction
  - Finds linear combinations of the features  $X = \{x_1, \dots, x_d\}$  with large ratios of between-groups to within-groups sums of squares - **discriminant variables**;
- Classification
  - Predicts the class of an observation  $x$  by the class whose mean vector is closest to  $x$  in terms of the discriminant variables