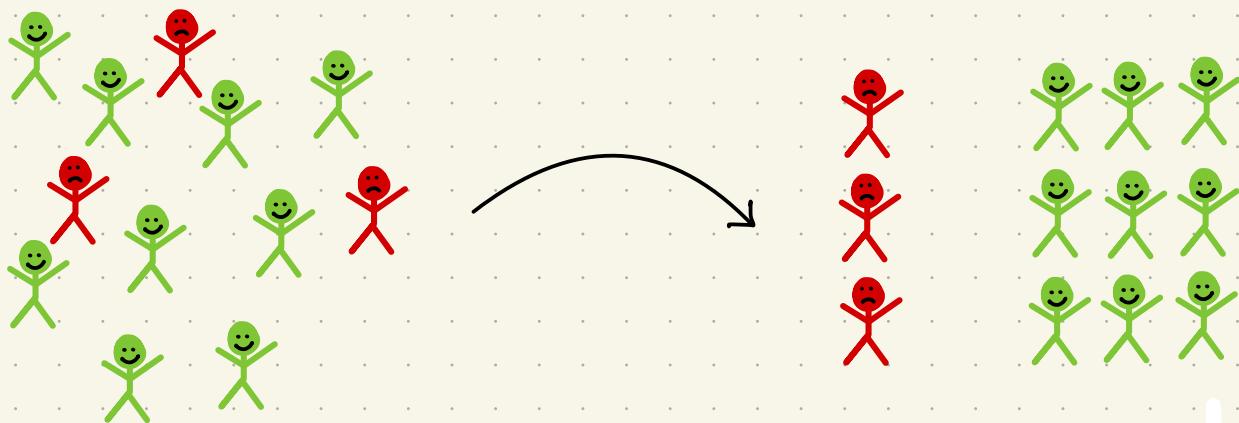


# Classification



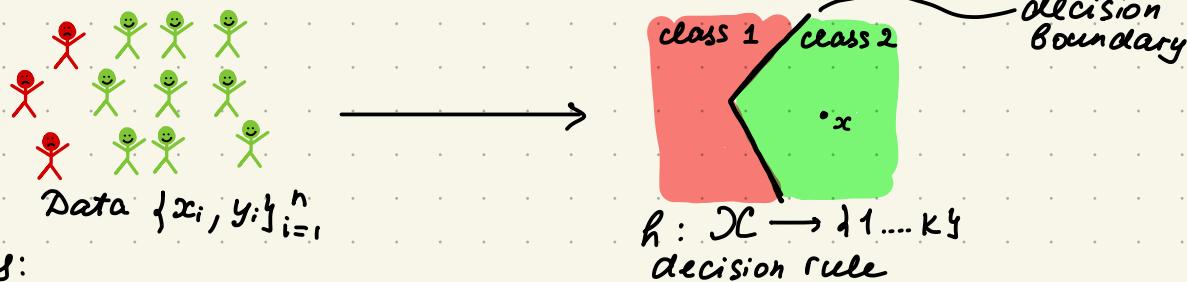
## Classification task

Given feature matrix  $X$  and response  $y$

$$X = \begin{pmatrix} -x_1^T- \\ \vdots \\ -x_n^T- \end{pmatrix} \in \mathbb{R}^{n \times p} \text{ and } y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

where  $y_i \in \{1, \dots, k\}$  gives class of  $i$ -th observation.

Denote by  $C_k \subseteq \{1, \dots, k\}$  the subset of observations that belong to class  $k$  and  $n_k = |C_k|$ .



### Examples:

- predict if a patient is sick ,  $y_i \in \{\text{sick}, \text{healthy}\}$
- predict image class  $y_i \in \{\text{cat}, \text{dog}, \text{bird}\}$
- predict movie genre  $y_i \in \{\text{comedy}, \text{thriller}, \text{drama}\}$

## Linear discriminant analysis

Denote by  $Z$  a random variable giving the class label.

$$Z = \begin{cases} 1 & \text{with probability } \pi_1 \\ \vdots \\ k & \text{with probability } \pi_k \end{cases} \quad \text{priors}$$

We want to build **Bayes classifier**

$$h(x) = \operatorname{argmax}_{k=1 \dots K} P(Z=k \mid X=x)$$

Equivalently,  $h(x) = \operatorname{argmax}_{k=1 \dots K} P(X=x \mid Z=k) \cdot \pi_k$

$$P(Z=k \mid X=x) = \frac{P(X=x \mid Z=k) \pi_k}{\sum_{j=1}^K P(X=x \mid Z=j) \pi_j} \leftarrow \text{common}$$

$$h(x) = \operatorname{argmax}_{k=1 \dots K} P(X=x \mid Z=k) \cdot \pi_k$$

Recall normal density:

$$f(x; \mu, \Sigma) = \left(\frac{1}{2\pi}\right)^{p/2} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Linear discriminant analysis (LDA) assumes

$$x \in \mathbb{R}^p \sim N_p(\mu_k, \Sigma) \quad \text{same variance}$$

- The decision rule is  $h(x) = \arg \max_{k=1 \dots K} (\alpha_k + b_k^T x)$
- $$h(x) = \arg \max_{k=1 \dots K} (\log f(x; \mu_k, \Sigma) + \log \pi_k) =$$
$$= \arg \max_{k=1 \dots K} \left( \log \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} - \frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log \pi_k \right)$$
$$= \arg \max_{k=1 \dots K} \left( -\frac{1}{2}(x^T \cancel{\Sigma^{-1}} x - 2\mu_k^T \Sigma^{-1} x + \mu_k^T \Sigma^{-1} \mu_k) + \log \pi_k \right) =$$
$$= \arg \max_{k=1 \dots K} (\alpha_k + b_k^T x) - \text{linear functions of } x$$

with  $\alpha_k = \log \pi_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k$  and  $b_k = \mu_k^T \Sigma^{-1}$

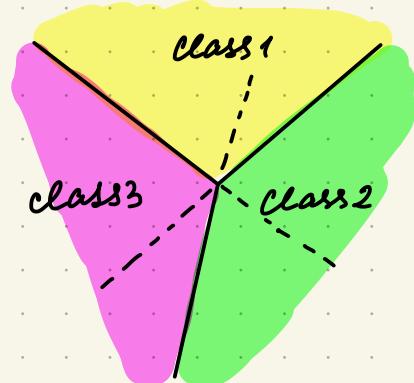
- The decision boundaries are linear.

Let's take two classes  $k, k'$

$$i \in C_k \text{ if } (\alpha_k + b_k^T x) > (\alpha_{k'} + b_{k'}^T x)$$

$$\text{Thus, } \alpha_k - \alpha_{k'} + (b_k - b_{k'})^T x > 0$$

(it's a half-space)



- In practice, we don't know  $\pi_k, \mu_k, \Sigma$ , so we need to estimate them.

$$\hat{\pi}_k = \frac{n_k}{n} \quad \hat{\mu}_k = \frac{1}{n_k} \sum_{i \in C_k} x_i$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^K \sum_{i \in C_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T = \sum_{k=1}^K \frac{n_k}{n} S_k$$

where  $S_k = \frac{1}{n_k} \sum_{i \in C_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$  - sample covariance for  $C_k$

- If  $\Sigma$  is not common for groups, we apply quadratic discriminant analysis (QDA)

$$h(x) = \underset{k=1..K}{\operatorname{argmax}} \left( -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k \right)$$

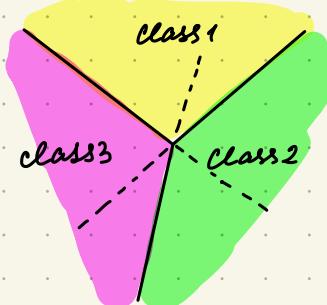
$$-\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k =$$

$$-\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_{k'})^T \Sigma_{k'}^{-1} (x - \mu_{k'}) + \log \pi_{k'}$$

$$\begin{aligned} & -\frac{1}{2} x^T (\Sigma_k^{-1} - \Sigma_{k'}^{-1}) x + \underbrace{(\mu_k^T \Sigma_k^{-1} - \mu_{k'}^T \Sigma_{k'}^{-1})}_P x - \frac{1}{2} (\mu_k^T \Sigma_k^{-1} \mu_k - \mu_{k'}^T \Sigma_{k'}^{-1} \mu_{k'}) \\ & + \log \pi_k - \log \pi_{k'} = x^T Q x + P^T x + R - \text{quadratic function} \\ & \quad R \end{aligned}$$

Estimates for  $\hat{\Sigma}_k = S_k = \frac{1}{n_k} \sum_{i \in C_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$

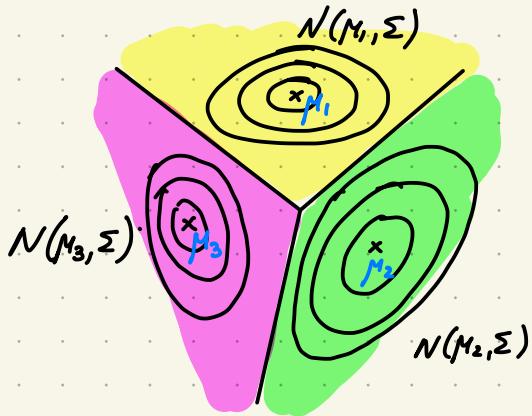
LDA



QDA



## Geometry of LDA



Is the decision boundary the same as Voronoi tessellation?

$$h_v(x) = \operatorname{argmin}_{k=1 \dots K} \|x - \mu_k\|^2$$

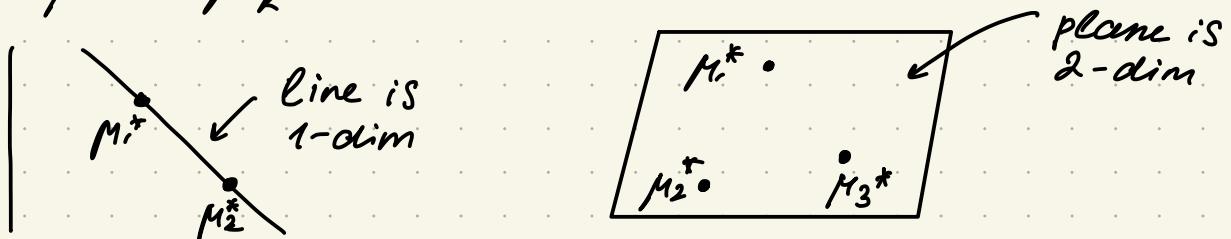
$$\begin{aligned} h_v(x) &= \operatorname{argmin}_{k=1 \dots K} \left( \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) - \log \pi_k \right) = \\ &= \operatorname{argmin}_{k=1 \dots K} \left( \frac{1}{2} (x - \mu_k)^T \Sigma^{-1/2} \Sigma^{-1/2} (x - \mu_k) - \log \pi_k \right) = \\ &= \operatorname{argmin}_{k=1 \dots K} \left( \frac{1}{2} \|\Sigma^{-1/2} (x - \mu_k)\|^2 - \log \pi_k \right) = \\ &= \operatorname{argmin}_{k=1 \dots K} \left( \frac{1}{2} \|x^* - \mu_k^*\|^2 - \log \pi_k \right) \end{aligned}$$

Where  $x^*$  and  $\mu_k^*$  are points in the transformed space

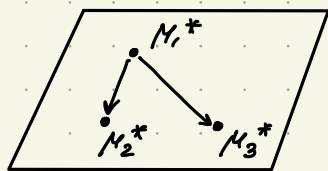
- The transformation  $x^* = \Sigma^{-1/2}x$  is spherizing,  
i.e.  $\text{cov}(x^*) = I$ .
- $\text{cov}(\Sigma^{-1/2}x) = \Sigma^{-1/2}\Sigma\Sigma^{-1/2} = I$
- Therefore, LDA classifies the points in the transformed space according to

$$h(x^*) = \arg \min_{k=1 \dots K} \left( \underbrace{\frac{1}{2} \|x^* - M_k^*\|^2}_{\text{nearest centroid}} - \underbrace{\log \pi_k}_{\text{adjustment by the class size.}} \right)$$

- $M_1^* \dots M_K^*$  lie in a plane  $M$  of dimension  $\leq K-1$ .



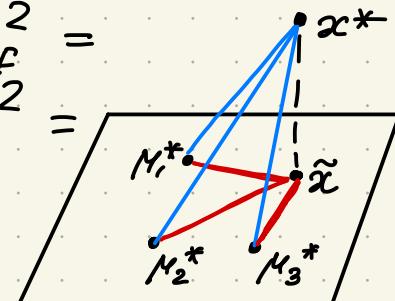
- Denote by  $P_M$  the projection operator onto this low-dimensional plane.



Find basis  $Q$  with  $Q^T Q = I$  for  $(M_2^* - M_1^*, \dots, M_k^* - M_1^*)$ . Then  $P_M = Q Q^T$ .

- For the projected points  $\tilde{x} = P_M x^*$  the value  $h(\tilde{x})$  is the same as  $h(x^*)$ .

$$\begin{aligned} \|x^* - M_k^*\|^2 &= \|P_M(x^* - M_k^*) + P_{M^\perp}(x^* - M_k^*)\|_F^2 = \\ &= \|P_M(x^* - M_k^*)\|^2 + \|P_{M^\perp}(x^* - M_k^*)\|^2 = \\ &= \|\tilde{x} - M_k^*\|^2 + \|P_{M^\perp}(x^*)\|^2 \end{aligned}$$



$$\begin{aligned} h(x^*) &= \underset{k=1 \dots K}{\operatorname{argmin}} \left( \frac{1}{2} \|x^* - M_k^*\|^2 - \log \pi_k \right) = \\ &= \underset{k=1 \dots K}{\operatorname{argmin}} \left( \frac{1}{2} \|\tilde{x} - M_k^*\|^2 + \frac{1}{2} \|P_{M^\perp}(x^*)\|^2 - \log \pi_k \right) = \\ &= \underset{k=1 \dots K}{\operatorname{argmin}} \left( \frac{1}{2} \|\tilde{x} - M_k^*\|^2 - \log \pi_k \right) = h(\tilde{x}) \end{aligned}$$

## LDA procedure:

1. Compute  $\hat{\pi}_k$ ,  $\hat{\mu}_k$ ,  $\hat{\Sigma}$
2. Sphere the data and transform centroids, project spherized data onto the plane  $M$  containing the spherized centroids.  
This can be combined in a single linear transformation

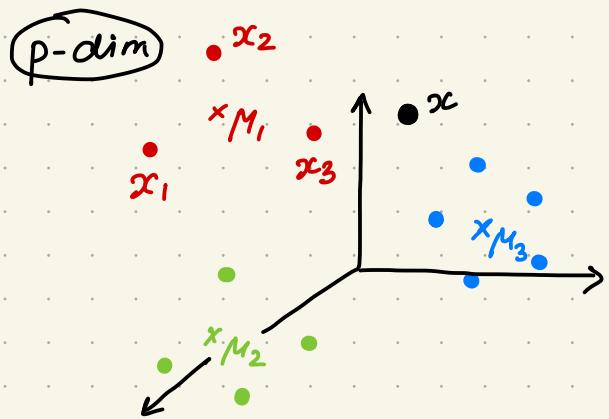
$$\tilde{x}_i = Ax_i \quad \tilde{\mu}_i = A\hat{\mu}_i$$

3. Given a new point  $x \in \mathbb{R}^P$  transform it to  $\tilde{x} = Ax$ , then classify to the nearest centroid  $\tilde{\mu}_i$ , adjusting for class proportions

- Decision boundaries in the  $M$  space are also linear.

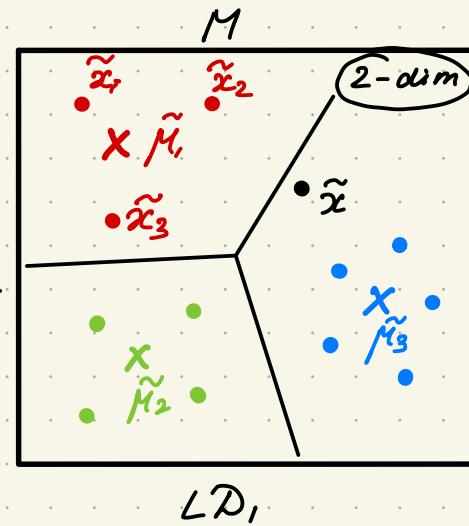
$$\begin{aligned}
 h(\tilde{x}) &= \underset{k=1 \dots K}{\operatorname{argmin}} \left( \frac{1}{2} \|\tilde{x} - \tilde{\mu}_k\|^2 - \log \pi_k \right) \\
 &= \underset{k=1 \dots n}{\operatorname{argmax}} \underbrace{\left( \log \pi_k - \frac{1}{2} \|\tilde{\mu}_k\|^2 \right)}_{\tilde{\alpha}_k} + \underbrace{\tilde{\mu}_k^T \tilde{x}}_{\tilde{b}_k} = \underset{k=1 \dots K}{\operatorname{argmax}} (\tilde{\alpha}_k + \tilde{\beta}_k^T \tilde{x})
 \end{aligned}$$

Example: If  $K=3$ ,  $M$  is two-dimensional.



$$A \in \mathbb{R}^{2 \times p}$$

~~~~~



## Reduced-rank DA (Fisher)

For  $K > 3$  we can find an  $L$ -dimensional subspace of  $M$  to project onto ( $L < K-1$ )  
 Choose the subspace that spread out the projected centroids (like in PCA!)

- $$S = \frac{1}{n} \sum_{k=1}^K \sum_{i \in C_k} (x_i - \bar{x})(x_i - \bar{x})^T =$$
 $\underbrace{\qquad\qquad\qquad}_{\text{total covariance}}$
- $$= \frac{1}{n} \underbrace{\sum_{k=1}^K \sum_{i \in C_k} (x_i - \bar{x}_k)(x_i - \bar{x}_k)^T}_{W - \text{within-class}} + \frac{1}{n} \underbrace{\sum_{k=1}^K \sum_{i \in C_k} n_k (\bar{x}_k - \bar{x})(\bar{x}_k - \bar{x})^T}_{B - \text{between-class}}$$

$$W = \sum_{k=1}^K \frac{n_k}{n} S_k \approx \sum_{k=1}^K \pi_k \Sigma_k (\Rightarrow \Sigma \text{ if common covariance})$$

$$B = \sum_{k=1}^K \frac{n_k}{n} (\bar{x}_k - \sum_{k=1}^K \frac{n_k}{n} \bar{x}_k) (\bar{x}_k - \sum_{k=1}^K \frac{n_k}{n} \bar{x}_k)^T \approx \sum_{k=1}^K \pi_k (\mu_k - \mu)(\mu_k - \mu)^T$$

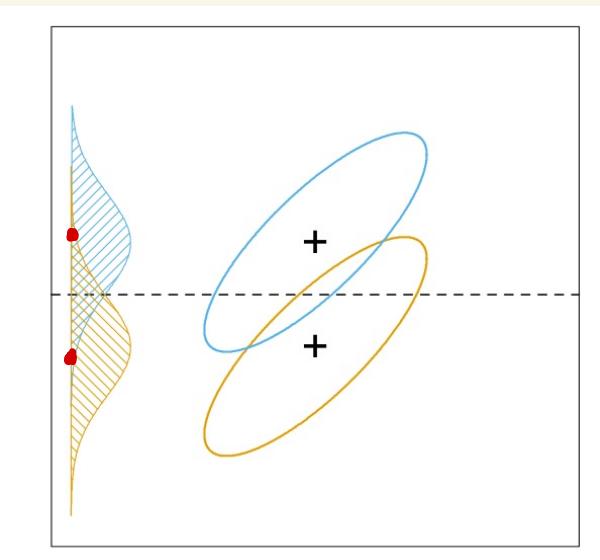
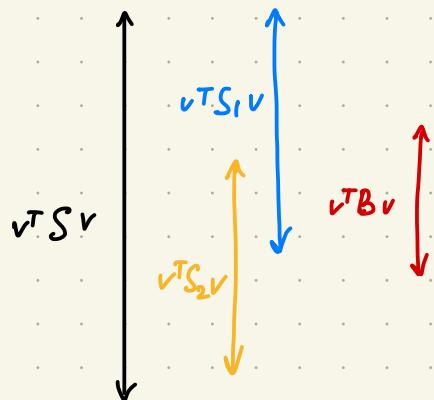
With  $\mu = \sum_{k=1}^K \pi_k \mu_k$  - Covariance & mean of centroids weighted by  $\pi_k$

- Suppose we want to project data onto direction  $v$ , then the variance of projections is

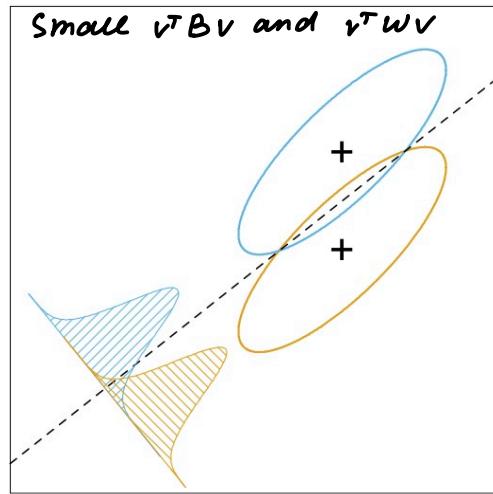
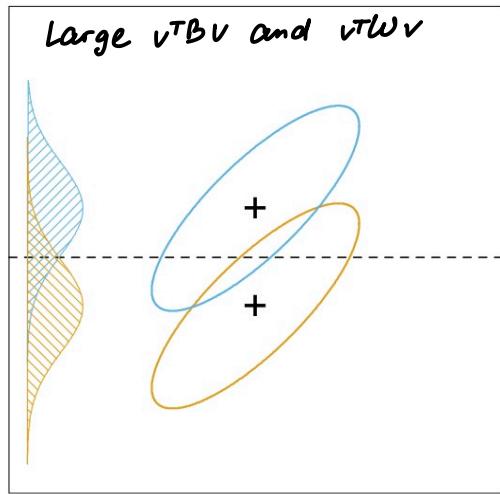
$$v^T S v = v^T W v + v^T B v$$

$$v^T W v = \sum_{k=1}^n \frac{n_k}{n} v^T S_k v = \sum_{k=1}^K \pi_k \cdot (\text{variance of projection in } C_k)$$

$v^T B v$  = variance of projections of centroids  
(weighted by  $\pi_k$ )



We want large  $v^T B v$  and small  $v^T W v$



Fisher's proposition:

maximize  
 $v \in R^p$

$$\frac{v^T B v}{v^T W v}$$

↑  
Rayleigh quotient

- The solution  $v$  is the largest eigenvector of  $W^{-1}B$ .

Let  $\tilde{v} = W^{1/2}v$  then  $v = W^{-1/2}\tilde{v}$  and Fisher's problem

$$\underset{\tilde{v} \in R^p}{\text{maximize}} \frac{\tilde{v}^T (W^{-1/2} B W^{-1/2}) \tilde{v}}{\|\tilde{v}\|^2}$$

Equivalently, if  $\tilde{B} = W^{-1/2} B W^{-1/2}$  we solve

$$\underset{\tilde{v} \in R^p}{\text{maximize}} \tilde{v}^T \tilde{B} \tilde{v} \text{ subject to } \|\tilde{v}\| = 1$$

Now,  $\tilde{v}$  is the largest e. vector of  $\tilde{B} \Rightarrow \tilde{B} \tilde{v} = \lambda, \tilde{v}$   
 ~~$W^{-1/2} B W^{-1/2} \tilde{v} = W^{-1/2} B \cancel{W^{1/2} W^{1/2}} v = W^{-1/2} B v = \lambda, W^{-1/2} v$~~

Then,  $W^{-1} B v = \lambda, v$ .

- If we assumed common covariance for classes ( $\Sigma \approx \Sigma$ ), then  $\tilde{B} = \sum_{k=1}^K \pi_k (\tilde{\mu}_k - \tilde{\mu})(\tilde{\mu}_k - \tilde{\mu})^T$  with  $\tilde{\mu}_k = A \mu_k$  where  $A$  is spherling + projection operator from LDA.

$$\begin{aligned} \Sigma^{-1/2} \mu_k &= \mu_k^* = \tilde{\mu}_k = A \mu_k \\ \tilde{B} &= \Sigma^{-1/2} B \Sigma^{-1/2} = \sum_{k=1}^K \pi_k \underbrace{\Sigma^{-1/2}(\mu_k - \mu)}_{\tilde{\mu}_k} \underbrace{(\mu_k - \mu)^T}_{\tilde{\mu}_k^T} \Sigma^{-1/2} = \\ &= \sum_{k=1}^K \pi_k (\tilde{\mu}_k - \tilde{\mu})(\tilde{\mu}_k - \tilde{\mu})^T \end{aligned}$$

- Thus  $\tilde{v}$  is the top eigenvector of  $\tilde{B}$  and is the first PC direction of  $\tilde{\mu}_1, \dots, \tilde{\mu}_K$  (where observations in PCA are weighted by  $\pi_1, \dots, \pi_K$ )

## Summary (LDA+dimension reduction)

- Compute centroids  $M_1, \dots, M_k$
- Compute within class covariance  $W$
- Transform centroids  $\tilde{M}_k = W^{-1/2} M_k$
- Compute between class covariance  $\tilde{B}$
- Denote by  $\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_l$  the eigenvectors of  $\tilde{B}$
- Project data onto  $v_1 \dots v_e$  with  $v_e = W^{-1/2} \tilde{v}_e$   
 $\uparrow$   
LD directions