

Linear Discriminant Analysis

Instructor: Junghye Lee

Department of Industrial Engineering

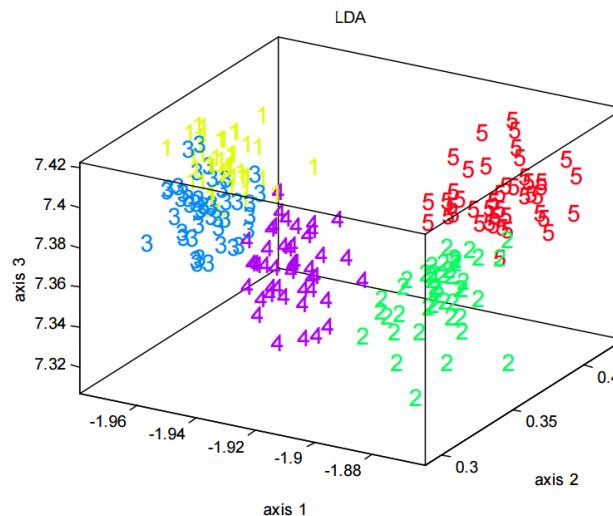
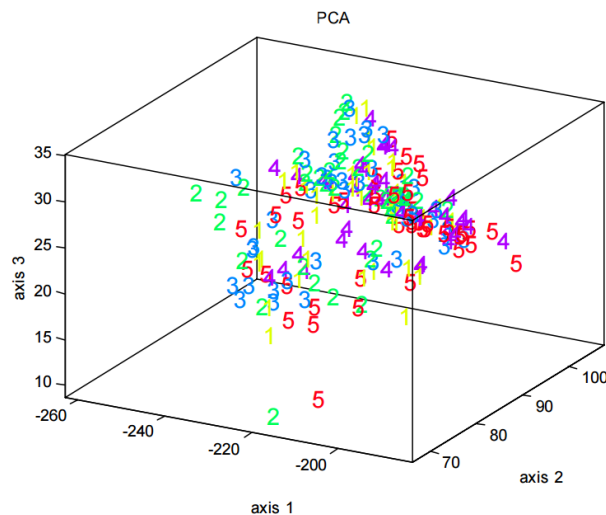
junghyelee@unist.ac.kr

Contents

- 1** Linear Discriminant Analysis, 2 classes
- 2** Linear Discriminant Analysis, C classes
- 3** LDA vs. PCA example
- 4** Limitations of LDA

Introduction

- Dimensionality reduction is a crucial concept in machine learning and data classification.
- The most famous example of dimensionality reduction is Principal Component Analysis (PCA):
 - Is an unsupervised method, so it doesn't include label information.
 - Searches for the directions the data have the largest variance.
 - There are difficulty issues with the number of principal components to choose.



Introduction

- Discriminant analysis methods can be good candidates to address such problems.
 - These methods are supervised, so they include label information.
 - The goal is to find directions on which the data is best separable.
- One of the very well-known discriminant analysis method is the Linear Discriminant Analysis (LDA).

Linear Discriminant Analysis - 2 classes

- The objective of LDA is to perform dimensionality reduction while preserving as much of the class discriminatory information as possible
 - Assume we have a D -dimensional N samples $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$, N_1 of which belong to class ω_1 , and N_2 to class ω_2
 - We seek to obtain a scalar vector \mathbf{y} by projecting the samples \mathbf{x} onto a line

$$y_j = \mathbf{w}^T \mathbf{x}_j \text{ where } j = 1, \dots, N$$

or

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

What is the dimension of \mathbf{x} , \mathbf{X} and \mathbf{w} ?

Illustration

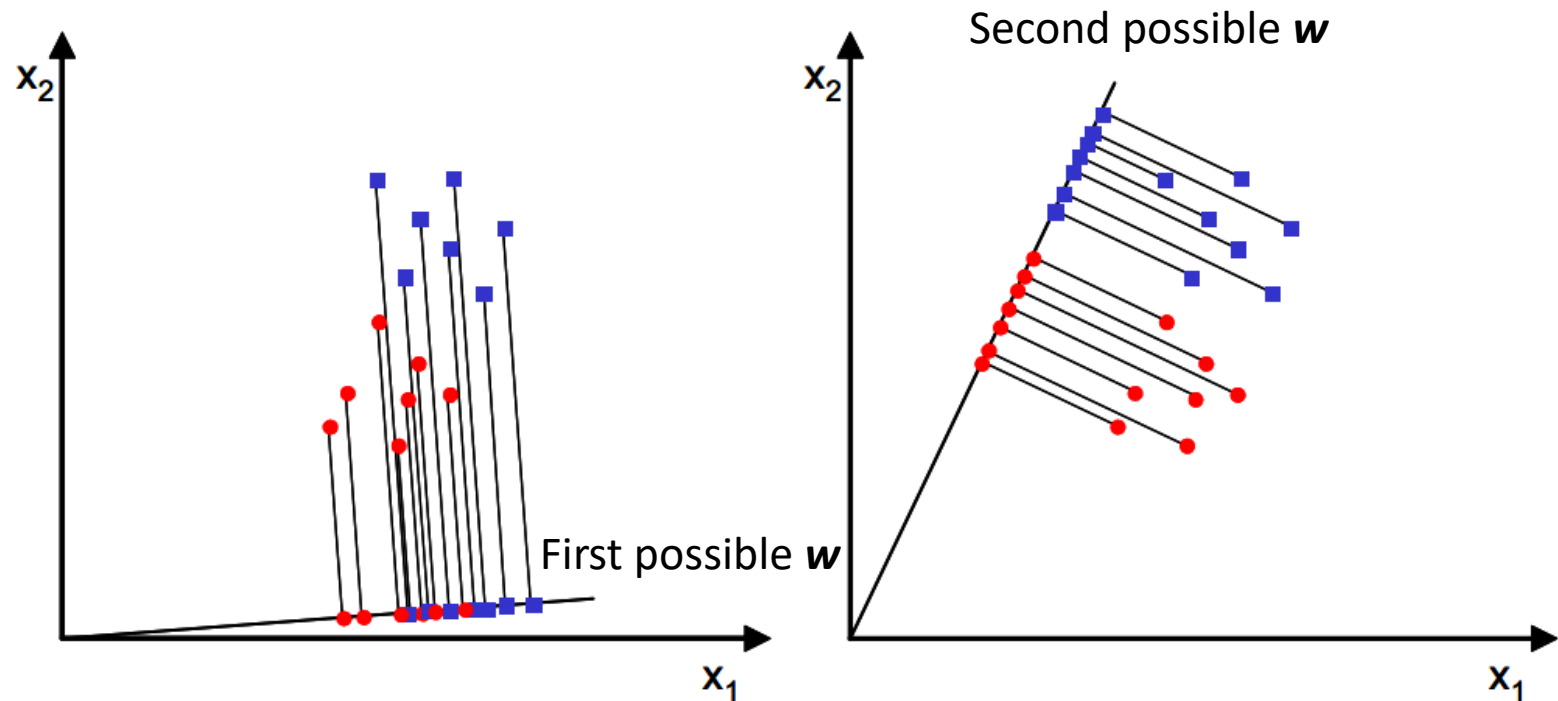
Sample #	AGE	gender	WAIST	BP_HIGH	BP_LWST	BLDS	Diabetes
1	44	1	86	120	80	75	0
2	56	0	84	110	70	151	0
3	38	1	78	103	61	82	1
4	60	1	88	130	77	153	1
5	28	1	92	128	77	101	0
6	42	1	94	134	85	91	0

parameters

$$\begin{aligned}
 \boxed{w^T} x_1 &= \text{Projection to } w \\
 (w_1, w_2, w_3, w_4, w_5, w_6) &\times \begin{pmatrix} 44 \\ 1 \\ 86 \\ 120 \\ 80 \\ 75 \end{pmatrix} = \boxed{w_1 \cdot 44 + w_2 \cdot 1 + w_3 \cdot 86 + w_4 \cdot 120 + w_5 \cdot 80 + w_6 \cdot 75} = y_1 \quad \text{scalar} \\
 \boxed{w^T} x_2 &= \\
 (w_1, w_2, w_3, w_4, w_5, w_6) &\times \begin{pmatrix} 56 \\ 0 \\ 84 \\ 110 \\ 70 \\ 151 \end{pmatrix} = \boxed{w_1 \cdot 56 + w_2 \cdot 0 + w_3 \cdot 84 + w_4 \cdot 110 + w_5 \cdot 70 + w_6 \cdot 151} = y_2 \\
 &\vdots
 \end{aligned}$$

Linear Discriminant Analysis - 2 classes

- Of all the possible lines we would like to select the one that maximizes the separability of the scalars
 - This is illustrated for the two-dimensional case in the following figures.



Which one is better?

Linear Discriminant Analysis – 2 classes

- In order to find a good projection vector, we need to define a measure of **separation between the projections**
 - The mean vector of each class in \mathbf{x} and \mathbf{y} feature space is

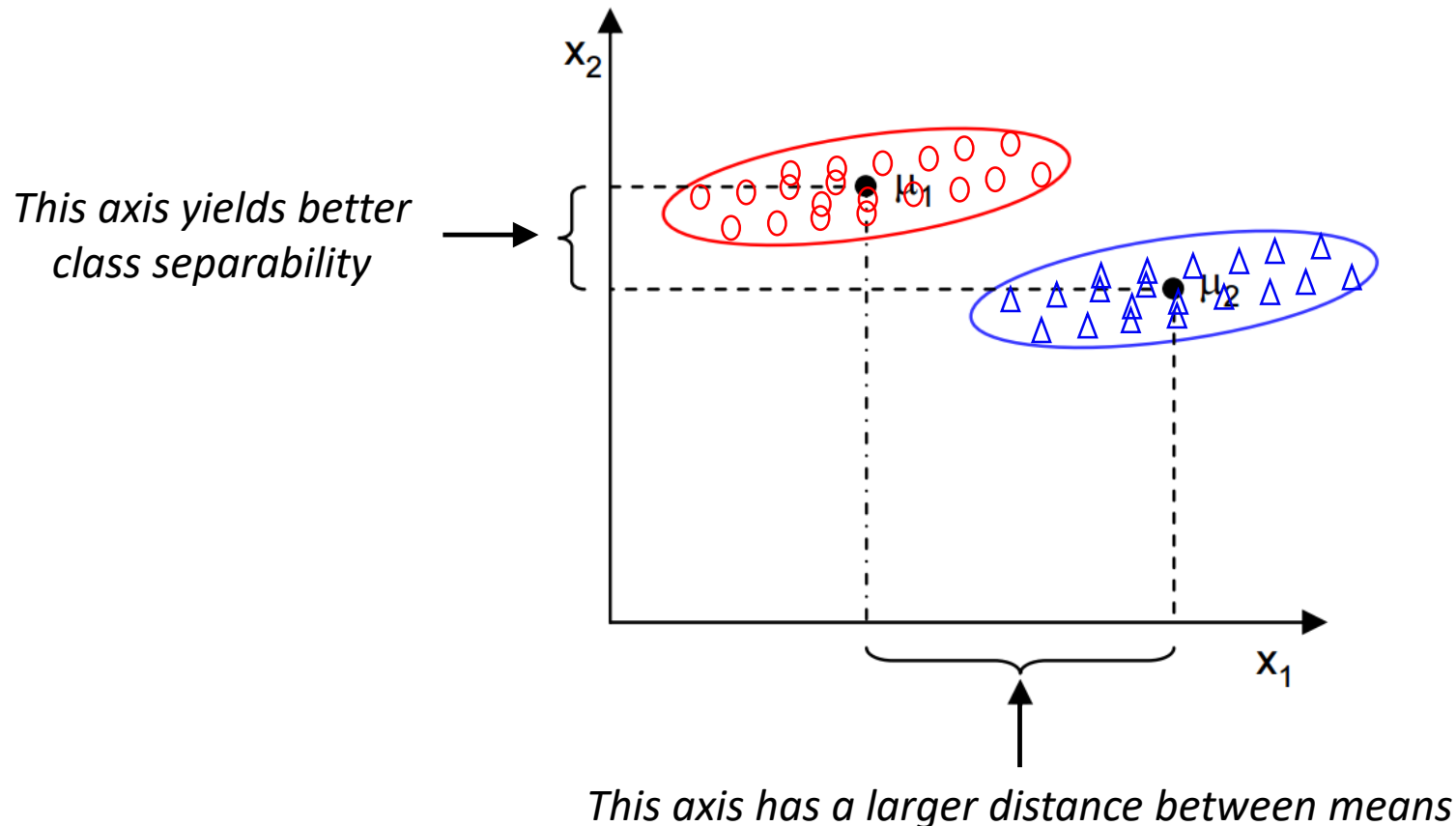
$$\boldsymbol{\mu}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in \omega_i} \mathbf{x} \text{ and } \tilde{\boldsymbol{\mu}}_i = \frac{1}{N_i} \sum_{\mathbf{y} \in \omega_i} \mathbf{y} = \frac{1}{N_i} \sum_{\mathbf{x} \in \omega_i} \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \boldsymbol{\mu}_i$$

- ω_i , where $i = 1, 2$. i.e., $\omega_1 = \text{class 1}$, $\omega_2 = \text{class 2}$
- We could then choose the distance between the projected means as our objective function

$$J(\mathbf{w}) = |\tilde{\boldsymbol{\mu}}_1 - \tilde{\boldsymbol{\mu}}_2| = |\mathbf{w}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)|$$

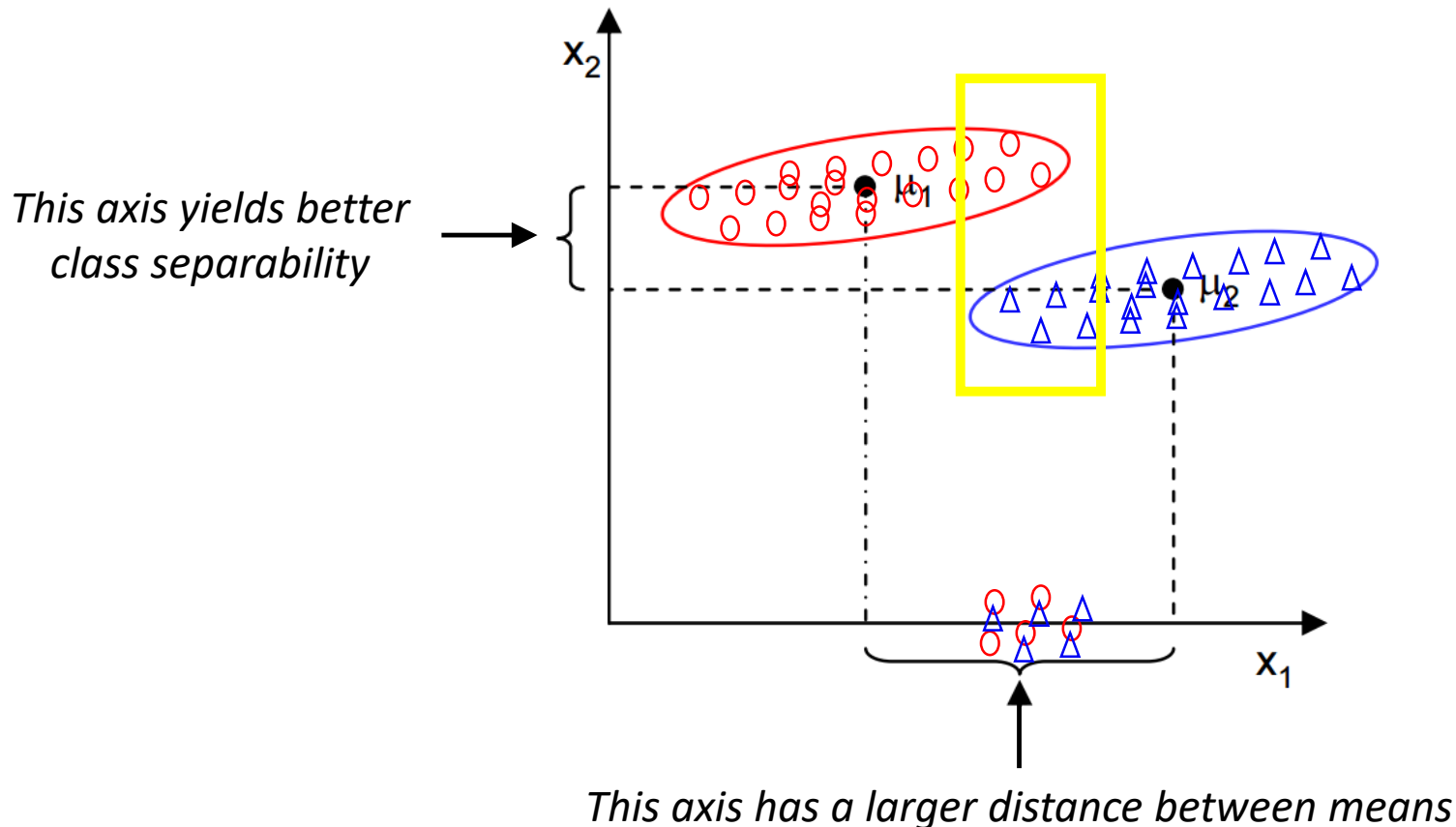
Linear Discriminant Analysis – 2 classes

- However, the distance between the projected means is not a very good measure since it does not take into account the standard deviation within the classes



Linear Discriminant Analysis – 2 classes

- However, the distance between the projected means is not a very good measure since it does not take into account the standard deviation within the classes



Linear Discriminant Analysis – 2 classes

- The solution proposed by Fisher is to maximize a function that represents the *difference between the means, normalized by a measure of the within-class scatter*
 - For each class we define the scatter, an equivalent of the variance, as

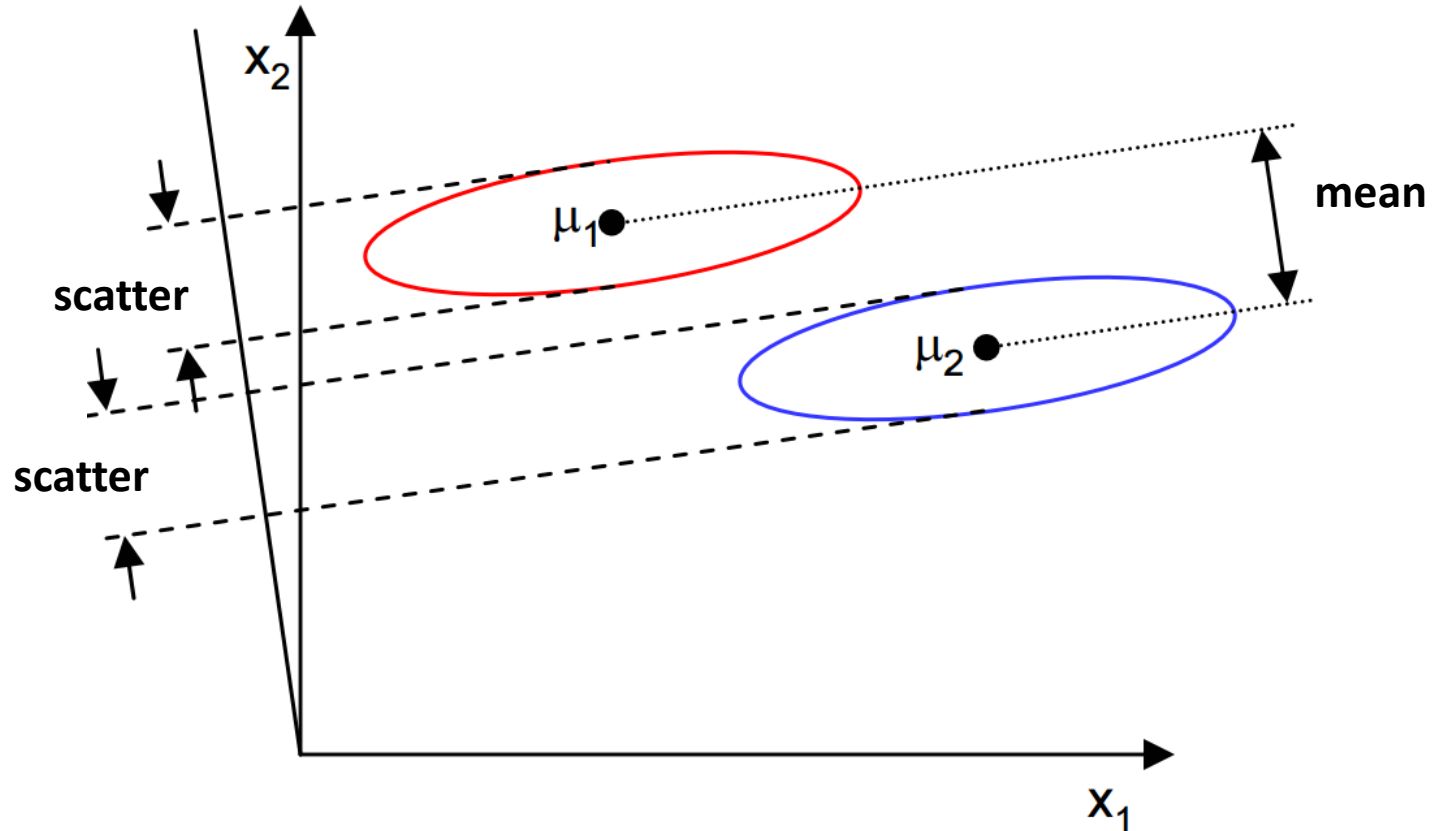
$$\tilde{s}_i^2 = \sum_{y \in \omega_i} (y - \tilde{\mu}_i)^2$$

- where the quantity $(\tilde{s}_1^2 + \tilde{s}_2^2)$ is called the within-class scatter of the projected examples
- The Fisher linear discriminant is defined as the linear function $\mathbf{w}^T \mathbf{x}$ that maximizes the criterion function

$$J(\mathbf{w}) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

Linear Discriminant Analysis – 2 classes

- Therefore, we will be **looking for a projection (w)** where **examples from the same class are projected very close to each other (*intra-class*)** and, at the same time, **the projected means between classes (*inter-class*)** are as far apart as possible.



Linear Discriminant Analysis – 2 classes

- In order to find the optimum projection \mathbf{w}^* , we need to express $J(\mathbf{w})$ as an explicit function of \mathbf{w}
- We define a measure of the scatter in multivariate feature space \mathbf{x} , which are scatter matrices

$$\mathbf{S}_i = \sum_{\mathbf{x} \in \omega_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T$$
$$\mathbf{S}_1 + \mathbf{S}_2 = \mathbf{S}_W$$

- where \mathbf{S}_W is called the within-class scatter matrix

Linear Discriminant Analysis – 2 classes

- The scatter of the projection \mathbf{y} can then be expressed as a function of the scatter matrix in feature space \mathbf{x}

$$\tilde{s}_i^2 = \sum_{y \in \omega_i} (y - \tilde{\mu}_i)^2 = \sum_{x \in \omega_i} (\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu}_i)^2 = \sum_{x \in \omega_i} \mathbf{w}^T (\mathbf{x} -$$

Linear Discriminant Analysis – 2 classes

- We can finally express the Fisher criterion in terms of S_W and S_B as

$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$$

- To find the maximum of $J(\mathbf{w})$ we derive and equate to zero

$$\begin{aligned}\frac{d}{d\mathbf{w}} [J(\mathbf{w})] &= \frac{d}{d\mathbf{w}} \left[\frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}} \right] = 0 \\ [\mathbf{w}^T S_W \mathbf{w}] \frac{d[\mathbf{w}^T S_B \mathbf{w}]}{d\mathbf{w}} - [\mathbf{w}^T S_B \mathbf{w}] \frac{d[\mathbf{w}^T S_W \mathbf{w}]}{d\mathbf{w}} &= 0 \\ [\mathbf{w}^T S_W \mathbf{w}] 2S_B \mathbf{w} - [\mathbf{w}^T S_B \mathbf{w}] 2S_W \mathbf{w} &= 0\end{aligned}$$

Linear Discriminant Analysis – 2 classes

Eigenvalue decomposition

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

$$\det(\mathbf{A} - \lambda\mathbf{I}) = 0$$

Generalized eigenvalue decomposition

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{B}\mathbf{v}$$

$$\det(\mathbf{A} - \lambda\mathbf{B}) = 0$$

- Dividing by $\mathbf{w}^T \mathbf{S}_W \mathbf{w}$

$$\left[\frac{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \right] \mathbf{S}_B \mathbf{w} - \left[\frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \right] \mathbf{S}_W \mathbf{w} = 0$$

$$\mathbf{S}_B \mathbf{w} - J \mathbf{S}_W \mathbf{w} = 0$$

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} - J \mathbf{w} = 0$$

to be full rank (invertible)

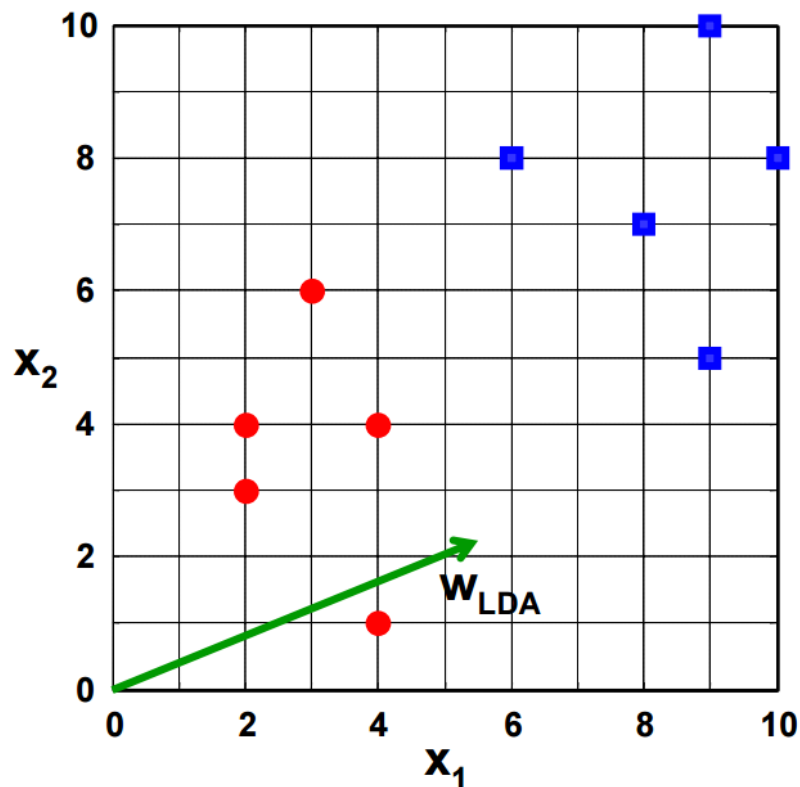
- Solving the **generalized eigenvalue problem** ($\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = J \mathbf{w}$) yields

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} \left\{ \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \right\} = \mathbf{S}_W^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

- This is known as Fisher's Linear Discriminant (1936), although it is not a discriminant but rather a specific choice of direction for the projection of the data down to one dimension

LDA example

- Compute the Linear Discriminant projection for the following two-dimensional dataset
 - $\omega_1: (x_1, x_2) = \{(4,1), (2,4), (2,3), (3,6), (4,4)\}$
 - $\omega_2: (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



LDA example

- Solution (by hand)
 - The class statistics are:

$$\mathbf{s}_1 = \begin{bmatrix} 4 & -2 \\ -2 & 13.2 \end{bmatrix}; \quad \mathbf{s}_2 = \begin{bmatrix} 9.2 & -0.2 \\ -0.2 & 13.2 \end{bmatrix}$$
$$\boldsymbol{\mu}_1 = [3.00 \quad 3.60]^T; \quad \boldsymbol{\mu}_2 = [8.40 \quad 7.60]^T$$

- The within- and between-class scatter are

$$\mathbf{S}_B = \begin{bmatrix} 29.16 & 21.60 \\ 21.60 & 16.00 \end{bmatrix}; \quad \mathbf{S}_W = \begin{bmatrix} 13.2 & -2.2 \\ -2.2 & 26.4 \end{bmatrix}$$

LDA example

- The LDA projection is then obtained as the solution of the generalized eigenvalue problem

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{v} = \lambda \mathbf{v}$$

$$|\mathbf{S}_W^{-1} \mathbf{S}_B - \lambda \mathbf{I}| = 0$$

$$\begin{bmatrix} 2.38 - \lambda & 1.76 \\ 1.02 & 0.75 - \lambda \end{bmatrix} = 0$$

$$\lambda = 3.13$$

$$\begin{bmatrix} 2.38 & 1.76 \\ 1.02 & 0.75 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 3.13 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0.91 \\ 0.39 \end{bmatrix}$$

- Or directly by

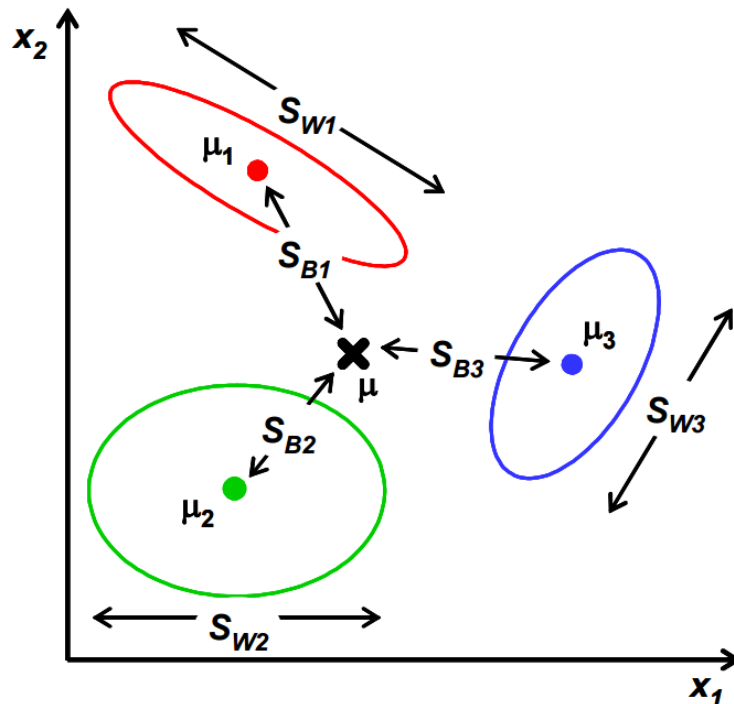
$$\mathbf{w}^* = \mathbf{S}_W^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = [0.91 \ 0.39]^T$$

Linear Discriminant Analysis – C classes

- Fisher's LDA generalizes very gracefully for C-class problems
 - Instead of one projection y , we will now seek $(C-1)$ projections $[y_1, y_2, \dots, y_{C-1}]$ by means of $(C-1)$ projection vectors w_i , which can be arranged by columns into a projection matrix

$$W = [w_1, w_2, \dots, w_{C-1}]:$$

$$y_i = w_i^T x \Rightarrow y = W^T x$$



Linear Discriminant Analysis – C classes

- Derivation

- The generalization of the within-class scatter is

$$S_W = \sum_{i=1}^C S_i$$

where $S_i = \sum_{x \in \omega_i} (x - \mu_i)(x - \mu_i)^T$ and $\mu_i = \frac{1}{N_i} \sum_{x \in \omega_i} x$

- The generalization of the total scatter is

$$S_T = \frac{1}{N} \sum_{\forall x} (x_i - \mu)(x_i - \mu)^T$$

where $\mu = \frac{1}{N} \sum_{\forall x} x = \frac{1}{N} \sum_i \sum_{x \in \omega_i} N_i \mu_i$

- The total scatter is

$$S_T = S_B + S_W$$

- The generalization of the within-class scatter is

$$S_B = \sum_{i=1}^C N_i (\mu_i - \mu)(\mu_i - \mu)^T$$

different from S_B for 2-class

Linear Discriminant Analysis – C classes

- Similarly, we define the mean vector and scatter matrices for the projected samples as

$$\begin{aligned}\tilde{\boldsymbol{\mu}}_i &= \frac{1}{N_i} \sum_{\mathbf{y} \in \omega_i} \mathbf{y} & \tilde{S}_W &= \sum_{i=1}^C \sum_{\mathbf{y} \in \omega_i} (\mathbf{y} - \tilde{\boldsymbol{\mu}}_i)^2 \\ \tilde{\boldsymbol{\mu}} &= \frac{1}{N} \sum_{\forall \mathbf{y}} \mathbf{y} & \tilde{S}_B &= \sum_{i=1}^C N_i (\tilde{\boldsymbol{\mu}}_i - \tilde{\boldsymbol{\mu}})(\tilde{\boldsymbol{\mu}}_i - \tilde{\boldsymbol{\mu}})^T\end{aligned}$$

- From our derivation for the two-class problem, we can write

$$\begin{aligned}\tilde{S}_W &= \mathbf{W}^T \mathbf{S}_W \mathbf{W} \\ \tilde{S}_B &= \mathbf{W}^T \mathbf{S}_B \mathbf{W}\end{aligned}$$

\mathbf{W} is a matrix!

What is the dimension?

Linear Discriminant Analysis – C classes

- Recall that we are looking for a projection that maximizes the ratio of between-class to within-class scatter. Since the projection is no longer a scalar (it has $C-1$ dimensions), we then use the determinant of the scatter matrices to obtain a scalar objective function:

$$J(W) = \frac{|\tilde{S}_B|}{|\tilde{S}_W|} = \frac{|W^T S_B W|}{|W^T S_W W|} = \text{trace} \left((\tilde{S}_W)^{-1} \tilde{S}_B \right)$$

- We will seek the projection matrix W^* that maximizes this ratio.
- It can be shown that the optimal projection matrix W^* is the one whose columns are the eigenvectors corresponding to the largest eigenvalues of the following generalized eigenvalue problem.

$$W^* = [w_1^*, w_2^*, \dots, w_{C-1}^*] = \operatorname{argmax}_W \left\{ \frac{|W^T S_B W|}{|W^T S_W W|} \right\} \Rightarrow (S_B - \lambda_i S_W) w_i^* = 0$$

(where $i = 1, \dots, C - 1$)

Linear Discriminant Analysis – C classes

- NOTES

- S_B is the sum of C matrices of **rank one or less (i.e., max C)** and the mean vectors are constrained by

$$\sum_{i=1}^C N_i (\boldsymbol{\mu}_i - \boldsymbol{\mu}) = \mathbf{0} \quad \text{We lose one rank.}$$

- Therefore, S_B will be of rank $(C-1)$ or less
 - This means that only $(C-1)$ of the eigenvalues λ_i will be non-zero
- The projections with maximum class separability information are the eigenvectors corresponding to the largest eigenvalues of $S_W^{-1} S_B$
- LDA can be derived as the Maximum Likelihood method for the case of **normal class-conditional densities with equal covariance matrices**

Decision Boundaries of LDA

- Okay, I understand how to find the projection for the C -class classification.
- Then, what is the decision boundary?
- Can you guess? Any ideas are welcome!

Decision Boundaries of LDA

- Once the transformation \mathbf{W} is given, the classification is then performed in the transformed space based on some distance metric, such as Euclidean distance and cosine measure:

$$\text{Euclidean distance } d(\mathbf{u}, \mathbf{v}) = \sqrt{\sum_j (u_j - v_j)^2}$$

$$\text{Cosine similarity } d(\mathbf{u}, \mathbf{v}) = 1 - \frac{\sum_j u_j v_j}{\sqrt{\sum_j u_j^2} \sqrt{\sum_j v_j^2}}$$

- Then upon the arrival of the new instance \mathbf{x}_{new} , it is classified to

$$\operatorname{argmin}_i = d(\mathbf{W}^T \mathbf{x}_{new}, \tilde{\boldsymbol{\mu}}_i) = d(\mathbf{W}^T \mathbf{x}_{new}, \mathbf{W}^T \boldsymbol{\mu}_i)$$

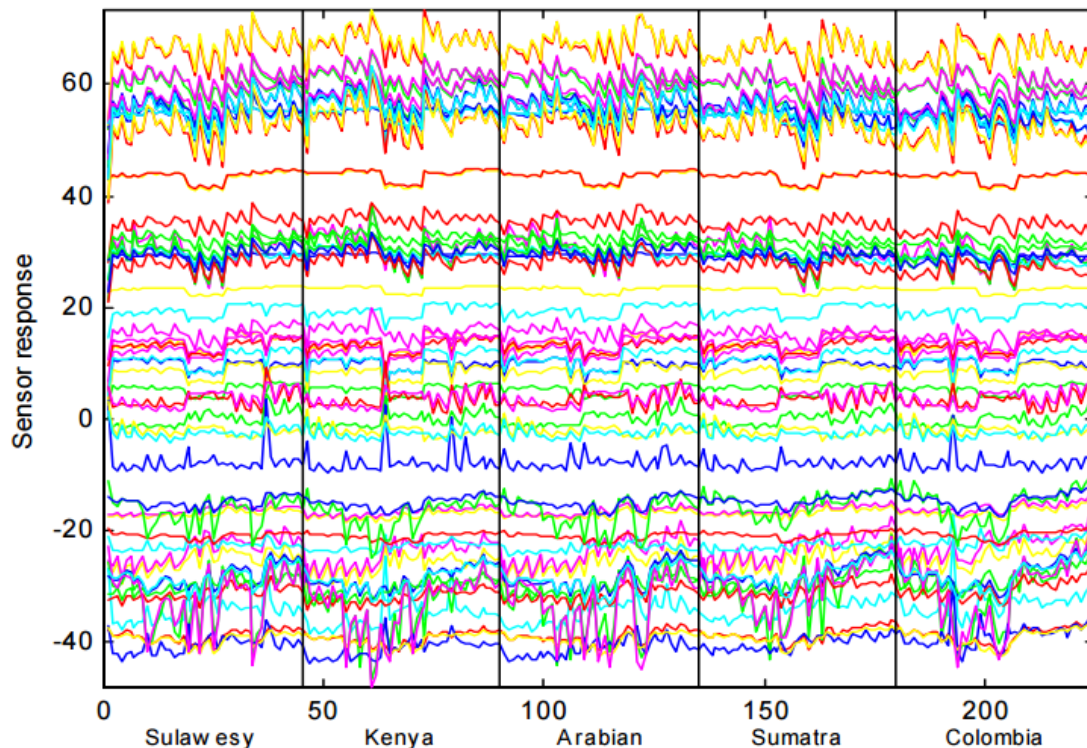
– where \mathbf{W} reduces to \mathbf{w} in binary classification.

Assumption of LDA

- LDA approaches the problem by assuming that the conditional probability density functions $p(x|y = 0)$ and $p(x|y = 1)$ are both normally distributed with mean and covariance parameters (μ_0, Σ_0) and (μ_1, Σ_1) , respectively.
- Regarding this, note that we can derive LDA by Bayes Theorem (please see the reference).

LDA vs. PCA Example

- Coffee discrimination with a gas sensor array
 - Five types of coffee beans were presented to an array of chemical gas sensors
 - For each coffee type, 45 “sniffs” were performed and the response of the gas sensor array was processed in order to obtain a 60-dimensional feature vector

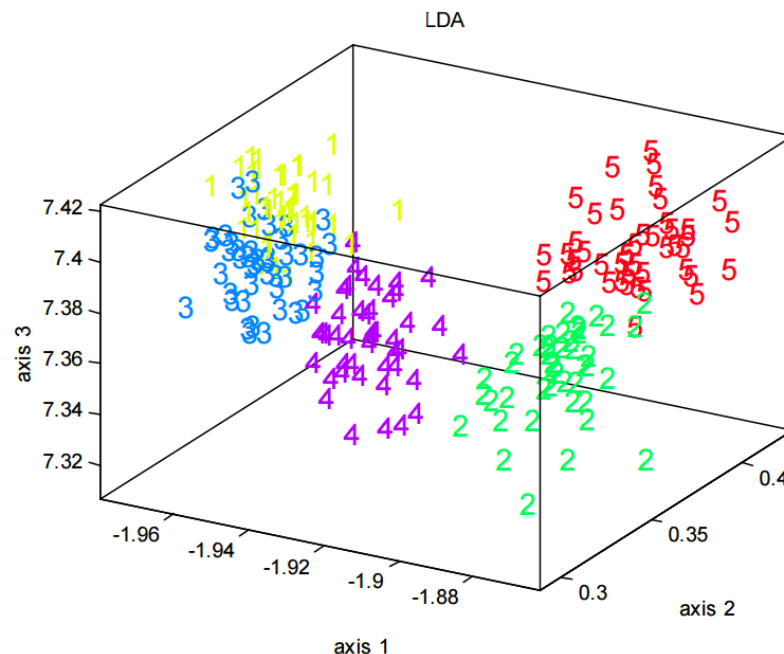
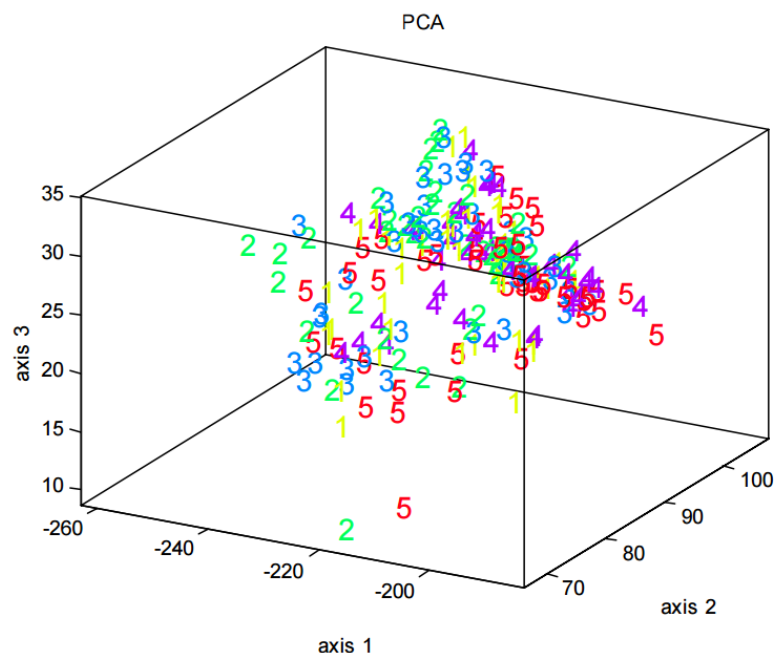


I.V. $X \in \mathbb{R}^{225 \times 60}$

D.V. $y_i \in \{1, 2, 3, 4, 5\}$
where $i = 1, \dots, 225$

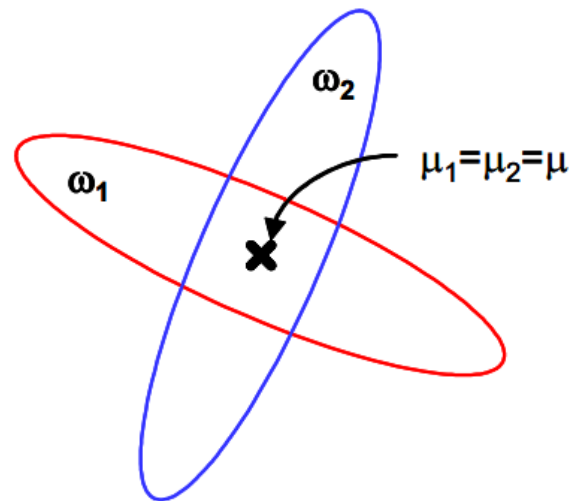
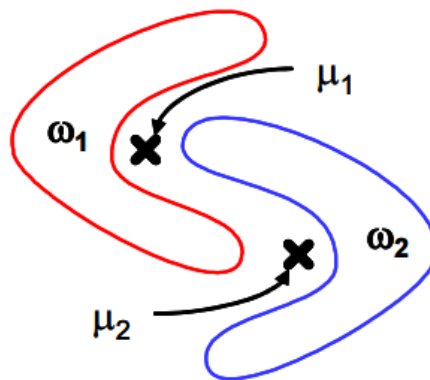
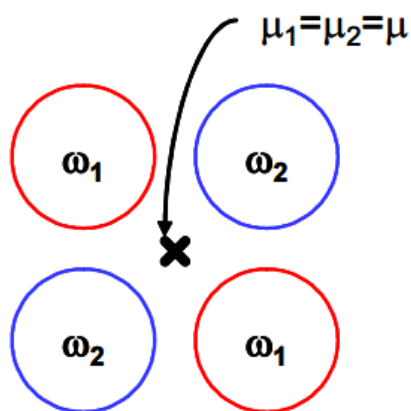
LDA vs. PCA Example

- Results
 - These figures show the performance of PCA and LDA on an **odor recognition problem**
 - From the 3D scatter plots **it is clear that LDA outperforms PCA** in terms of class discrimination
 - This is one example where the **discriminatory information is not aligned with the direction of maximum variance**



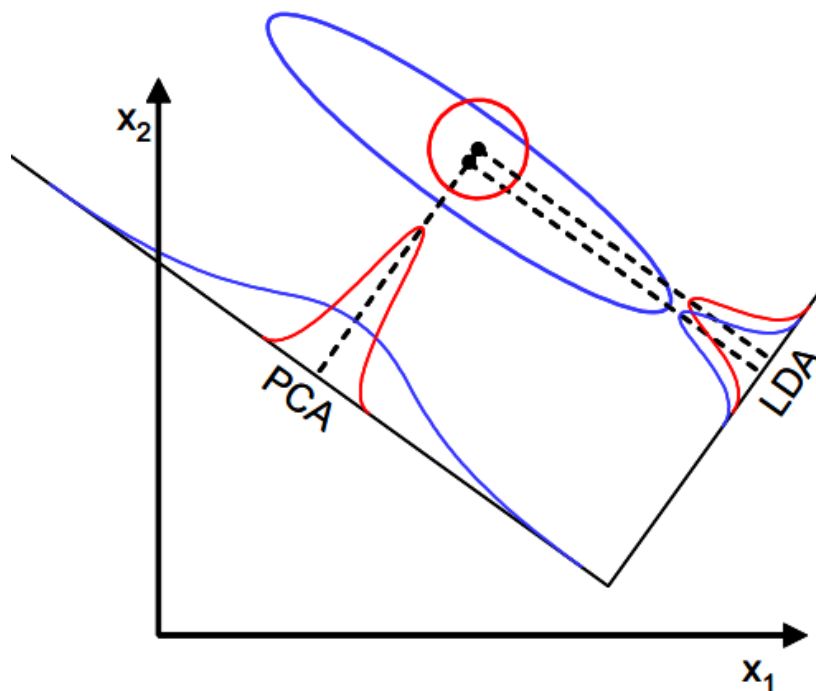
Limitations of LDA

- LDA produces at most $C-1$ feature projections
 - If the classification error estimates establish that more features are needed, some other method must be employed to provide those additional features
- LDA is a parametric method since it assumes unimodal Gaussian likelihoods
 - If the distributions are significantly non-Gaussian, the LDA projections will not be able to preserve any complex structure of the data, which may be needed for classification



Limitations of LDA

- LDA will fail when the discriminatory information is not in the mean but rather in the variance of the data.



**In this case, PCA is encouraged.
You should look at your data first.**

Questions?