

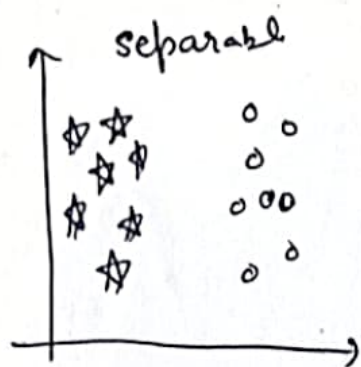
Principal Component Analysis Σ - matrix

- ① Calculate the covariance Σ of data point.
- ② Calculate eigenvectors + correspond eigenvalue.
- ③ Sort eigenvectors according to their given value in decreasing order.
- ④ Choose first k eigenvectors + that will be the new k dimension.
- ⑤ Transform original n -dim to k -dimension.

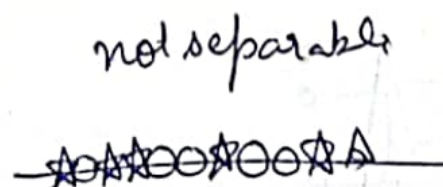
PCA is an unsupervised learning technique that aims to maximize the variance of the data along the principal components.

→ Aim is to identify the directions (components) that captures as much variance as possible.

However, the directions of maximum variance may be useless for classification.



PCA →



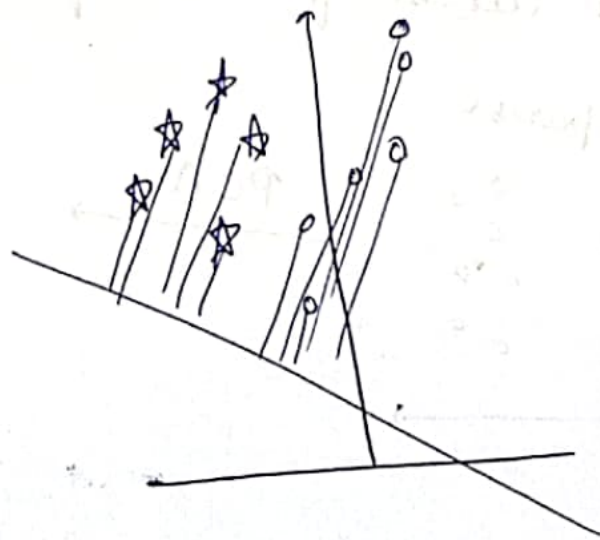
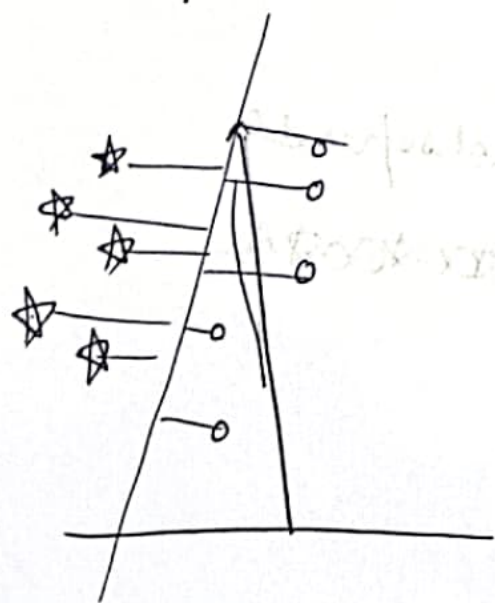
If the data is linearly separable in higher dimensional space then the linearly separable should be preserved in the lower dimensionality space.

Here LDA (Linear Discriminant Analysis) comes to the rescue.

→ Unlike ~~PCA~~ PCA, LDA is a supervised learning method, which means it takes class labels into account when finding direction of maximum variance.

LDA aims to maximize the separation b/w different classes in the data.

→ goal is to identify the direction that captures the most separation b/w the classes



Given a data set $x_1, x_2, \dots, x_n \in \mathbb{R}^d$ where

n_1 sample are coming from class C_1

n_2 — do ————— C_2

Our Aim is to find a unit vector direction that
"best discriminates" b/w classes

Consider any unit vector $v \in \mathbb{R}^d$

1D projection of points are

$$y_i = v^T x_i \quad i=1, 2, \dots, n$$

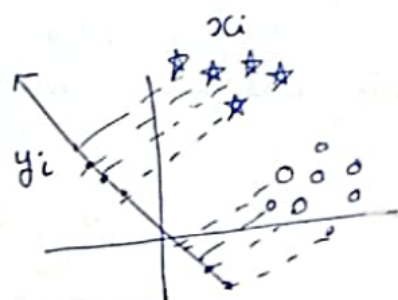
Let μ_1 & μ_2 be the mean of centroid of class
 C_1 and C_2 , respectively before projection

Let $\tilde{\mu}_i$ denote mean of sample of class C_i after the
projection then

$$\begin{aligned} \tilde{\mu}_1 &= \frac{1}{n_1} \sum_{x_i \in C_1}^{n_1} v^T x_i = v^T \left(\frac{1}{n_1} \sum_{x_i \in C_1}^{n_1} x_i \right) \\ &= v^T (\mu_1) \end{aligned}$$

$$\tilde{\mu}_2 = \frac{1}{n_2} \sum_{x_i \in C_2}^{n_2} v^T x_i = v^T (\mu_2)$$

What we solve in LDA is



maximize $|\tilde{\mu}_1 - \tilde{\mu}_2|$ where $\tilde{\mu}_j = \mu_j^T v$

$v: \|v\|=1$

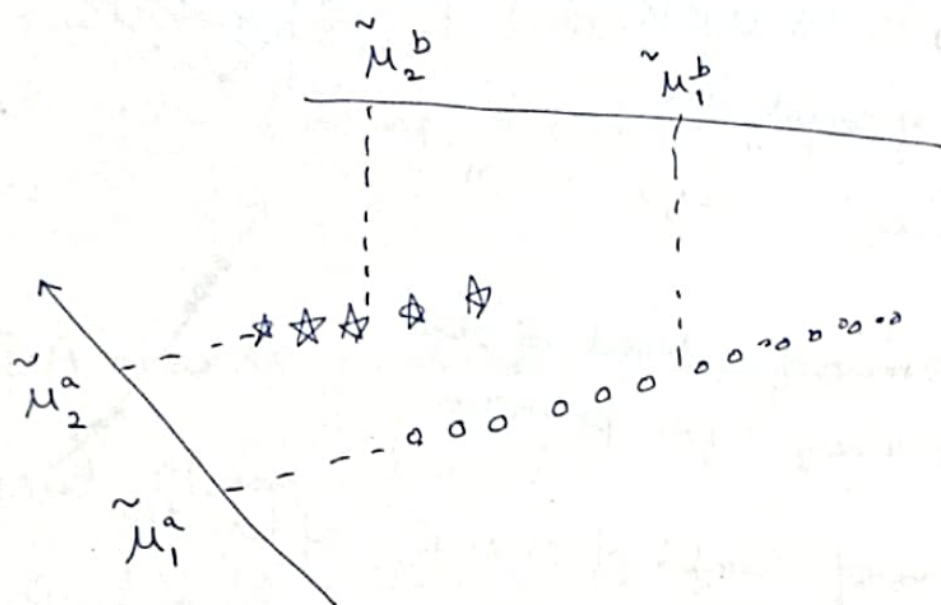
↓

$\tilde{\mu}_j = v^T \mu_j$
 $j=1,2$

Means

maximize $|\tilde{\mu}_1 - \tilde{\mu}_2|$ w.r.t v , while ensuring v remain unit vector

But this criteria might not always work



So we also need to pay attention to the variance of the projected class.

let $y_i = v^T x_i$ be the projected samples.

Then scatter for sample of class C_1 is

$$\tilde{S}_1^2 = \sum_{x_i \in C_1} (y_i - \tilde{\mu}_1)^2 \quad \text{lets ignore } \frac{1}{n} \text{ for now}$$

$$\tilde{S}_2^2 = \sum_{x_i \in C_2} (y_i - \tilde{\mu}_2)^2$$

We use them to normalize $(\tilde{\mu}_1, \tilde{\mu}_2)$ by variance.

Ideally, projected classes have both

- far away means
- small variances.

Thus, one need to project the data onto a line having direction v which maximizes

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

class 1 and class 2 scatter after project should be small

$$\max_{v: \|v\|=1} \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{s}_1^2 + \tilde{s}_2^2} \quad \begin{array}{l} \text{optimal } v \text{ should.} \\ \rightarrow (\tilde{\mu}_1 - \tilde{\mu}_2) \text{ large} \\ \rightarrow (\tilde{s}_1^2 + \tilde{s}_2^2) \text{ small} \end{array}$$

$$\begin{aligned} (\tilde{\mu}_1 - \tilde{\mu}_2)^2 &= (v^T \mu_1 - v^T \mu_2)^2 = (v^T (\mu_1 - \mu_2))^2 \\ &= v^T (\mu_1 - \mu_2) \cdot (\mu_1 - \mu_2)^T v \\ &= v^T S_B v \end{aligned}$$

where $S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T \in \mathbb{R}^{d \times d}$ ~~is called~~
↓
called between-class scatter η_B

Remark: S_B is square, symmetric and positive semidefinite η_B .

Moreover $\text{rank}(S_B) = 1, \Rightarrow$ only 1 true eigenvalue.

Define separate class scatter S_1 & S_2 of class C_1 & C_2

$$S_1 = \sum_{x_i \in C_1} (x_i - \mu_1)(x_i - \mu_1)^T$$

$$S_2 = \sum_{x_i \in C_2} (x_i - \mu_2)(x_i - \mu_2)^T$$

Again $\frac{1}{n}$ is
not taken.

Now define within class scatter S_W as

$$S_W = S_1 + S_2$$

For each class $j=1,2$, variance of projection (onto u) is

$$\tilde{S}_j^2 = \sum_{x_i \in C_j} (y_i - \tilde{\mu}_j)^2 = \sum_{x_i \in C_j} (u^T x_i - u^T \mu_j)^2$$

$$= \sum_{x_i \in C_j} u^T (x_i - \mu_j)(x_i - \mu_j)^T u$$

$$\tilde{S}_1^2 = \sum_{x_i \in C_1} u^T (x_i - \mu_1)(x_i - \mu_1)^T u$$

$$\tilde{S}_2^2 = \sum_{x_i \in C_2} u^T (x_i - \mu_2)(x_i - \mu_2)^T u$$

$$\tilde{S}_j^2 = u^T \left[\sum_{x_i \in C_j} (x_i - \mu_j)(x_i - \mu_j)^T \right] u$$

$$= u^T S_j u$$

$$\tilde{S}_1^2 = u^T S_1 u$$

$$\tilde{S}_2^2 = u^T S_2 u$$

$$\begin{aligned}\tilde{S}_1^2 + \tilde{S}_2^2 &= v^T S_1 v + v^T S_2 v \\ &= v^T (S_1 + S_2) v \\ &= v^T S_w v\end{aligned}$$

$$S_w = S_1 + S_2 =$$

$$\sum_{x_i \in C_1} (x_i - \mu_1)(x_i - \mu_1)^T + \sum_{x_i \in C_2} (x_i - \mu_2)(x_i - \mu_2)^T$$

S_w is called total within-class scatter of original data

$S_w \in \mathbb{R}^{d \times d}$ is also square, symmetric & positive semidefinite

Putting everything together

$$\max_{v: \|v\|=1} \frac{v^T S_B v}{v^T S_w v} = J(v)$$

$$\frac{d}{d(v)} J(v) = 0 \quad \text{after some step}$$

$$\Rightarrow S_B v - \frac{v^T S_B v}{v^T S_w v} (S_w v) = 0$$

scalar value.
assume it to be λ

$$S_B v - \lambda S_w v = 0$$

$$S_B v = \lambda S_w v$$

$$(S_w^{-1} S_B) v = \lambda v$$

Matrix, let say M

$$M v = \lambda v$$

→ eigenvector
& eigenvalue
largest eigenvalue

v is the eigenvector of $S_w^{-1} S_B$ corresponding to largest eigenvalue.

However it is not computationally efficient method as .

- ① First one has to invert $S_w \rightarrow S_w^{-1}$
- ② Multiply it to $S_B \rightarrow S_w^{-1} S_B$
- ③ & then solve for eigenvalue.

$$S_w^{-1} S_B v = \lambda v$$

One can do it in a smarter way

$$S_w^{-1} S_B v = \lambda v$$

$$\lambda v = S_w^{-1} S_B v$$

we know $S_B x$ points in the same direction as

$$\mu_1 - \mu_2$$

$$S_B x = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T x$$

$$\lambda v = S_w^{-1} (\mu_1 - \mu_2) \underbrace{(\mu_1 - \mu_2)^T v}_{\text{scalar}}$$

This implies that $v \propto S_w^{-1} (\mu_1 - \mu_2)$

and it can be computed from $S_w^{-1} (\mu_1 - \mu_2)$

through rescaling! No need to calculate S_B

Recipe for LDA

- ① Calculate the mean vector for each class.

Mean vector class 1 $\rightarrow \mu_1$ (μ_1^x, μ_1^y)

Mean vector class 2 $\rightarrow \mu_2$ (μ_2^x, μ_2^y)

- ② Calculate the within-class scatter η

$$S_w = \sum_{x_i \in C} (x_i - \mu_c)(x_i - \mu_c)^T \text{ for all data points } x_i \text{ in class } C \text{ (} c = 1, 2 \text{)}$$

- ③ Calculate between-class scatter η .

$$S_b = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T$$

- ④ Compute eigenvalues & eigenvectors of η .

$$S_w^{-1} \cdot S_b$$

Eigenvalues represent separability of classes, and the corresponding eigenvectors are direction in which data should be projected

- ⑤ Sort top k eigenvectors

- ⑥ Create a $n \times k$ W with selected eigenvectors as columns

Multiple Discriminal Analysis (MDA)

Just like LDA, MDA aims to maximize the ratio of b/w-class scatter (S_B) to within-class scatter (S_W).

However, instead of finding a single vector v , it seeks a matrix V whose columns are discriminant function weights.

Soln is found by solving

$$S_W^{-1} S_B V = \Lambda V$$

columns of V are eigenvectors, which represent discriminant functions.

Λ - diagonal matrix of eigenvalues.

	PCA	LDA
Label	unsupervised	Supervised
Criterion	Variance	discrimination
classifier	feature classification	data classification
dimension	reduce data to upto any dimension	reduce data to upto "number of class - 1" dimension
linear projection	Yes	Yes