# Linear Discriminent Analysis (LDA)

Linear model for classification and dimentionality reduction

Useful for feature extraction in pattern recognition face recognition

LPA is a stastical technique for catagorizing data into groups.

By maximizing the seperation between classes, it enables accurate classification of new data points.

Logistic regression falls short in multiclass classification where LDA shines.
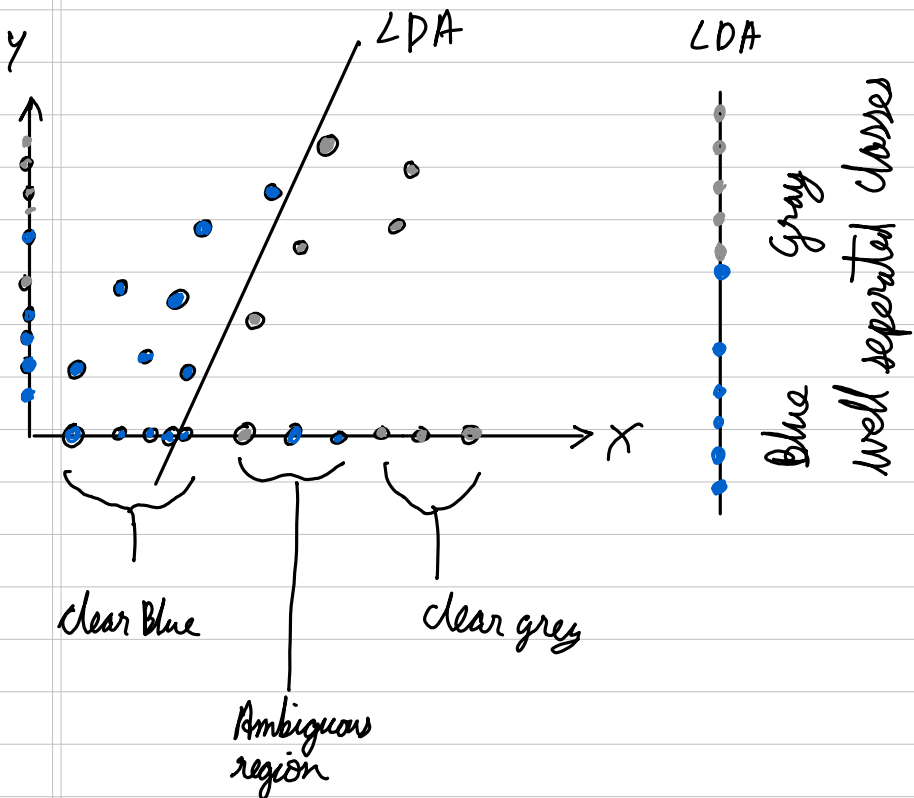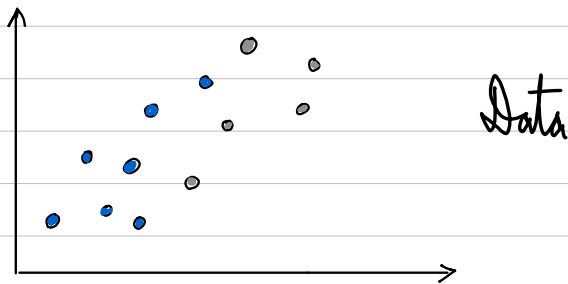
# What LDA does

LDA is just like PCA with a different objective

In PCA, focus is on the line that captures the most variation.

In LDA, we have labelled data. We know which class of the data it belongs to. LDA is supervised while PCA is unsupervised

The aim of LDA is to maximize the seperability between 2 classes (known)

We know beforehand which point is in which class unlike t-SNE



Data



Y

LDA

LDA

Blue Gray

well seperated classes

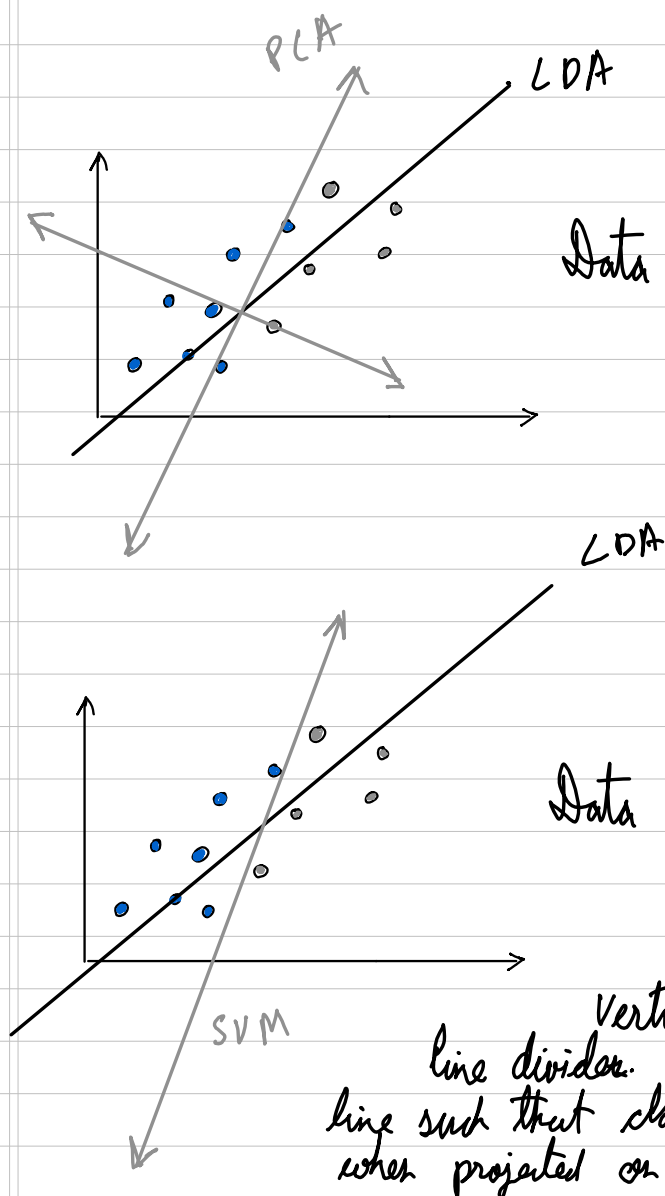clear Blue

clear grey

Ambiguous region

LDA uses information of the classes to
create a new axis that maximizes the separation
between the two classes.

like PCA & Unlike t-SNE, LDA
just chooses a line to project on. It
does not do any other change in the
points.

t-SNE on the otherhand shuffles the
points & then tries to match with the
original ordering.

Goal is that even after dimentionality reduction
the classes must remain classes, but
when classes are known.



SVM splits data
Vertically accross as a
line divider. But LDA draws
line such that classes get divided
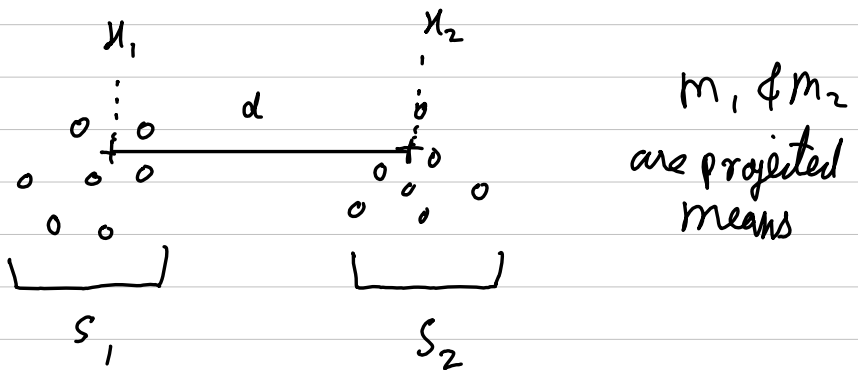when projected on to a lower dimention

# Well seperated classes

what does well seperated classes mean?

It means they must be far apart ①
and elements in class must be close
together ②

① Maximize the distance between 2
class means

② Minimize the variation (scatter) in between
a class



$m_1$ & $m_2$
are projected
means

$m_1 - m_2 \longrightarrow$ large

$S_1 ; S_2 \longrightarrow$ small
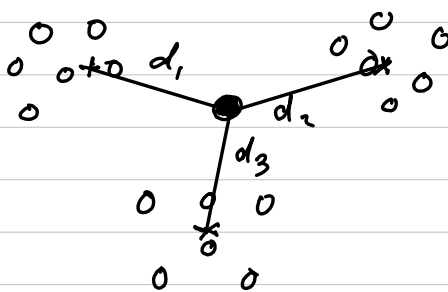
Hence the objective becomes

Maximize $\dfrac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$

| This is known as Fisher Discriminant ratio

Note → If we dont consider condition ②
and maximize the distance between the
means without considering the scatter, then
the seperation between the classes is not that
great! The classes will be fuzzy. overlapping

for Multiple classes, first find the
centroid of the whole dataset. Then
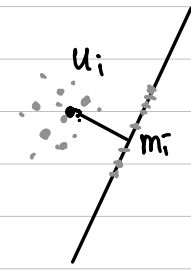maximize the distance of mean of
every class from the datapoint



Maximize $\dfrac{d_1^2 + d_2^2 + d_3^2}{S_1^2 + S_2^2 + S_3^2}$

$S \rightarrow$ Variance of projection

$x \rightarrow$ initial centroid of all data
$m \rightarrow$ projected centroid of all data

$m_i \rightarrow$ Projected Mean of class $i$ (Projected Centroid)
$x_i \rightarrow$ Initial mean of class $i$ (Initial Centroid)

$$y = Wx \rightarrow \text{Projection line}$$

$$\therefore \; m_i = W^T \mu_i$$



### Numerator of objective $\rightarrow$ Maximize

$$d_i^2 = (m_i - m)^2 = (W^T \mu_i - W^T \mu)^2$$

(Inter class scatter projected)
$$= W^T (\mu_i - \mu)(\mu_i - \mu)^T W$$

Inter class scatter Initial $\rightarrow S_{b_i}$

$$S_{b_i} = (\mu_i - \mu)(\mu_i - \mu)^T$$

$$\therefore \; d_i^2 = W^T S_{b_i} W$$

$$d_1^2 + d_2^2 + \cdots = W^T S_B W$$
$\uparrow$
Matrix of all classes

### Denominator $\rightarrow$

Within class scatter matrix Initial $\rightarrow$

for class $i$ $\quad S_i' = \sum_{x_i \in class i} (x_i - \mu_i)(x_i - \mu_i)^T$

Variance for each class (projected) $\rightarrow$

for class $i$ $\quad S_i^2 = \sum_{x_i \in class i} (W^T x_i - m_i)^2$

$S_i^2$ can be re written as $\rightarrow$

$$\tilde{S_i^2} = \sum_{x_i \in \mathcal{U}_i} (W^T x_i - W^T \mu_i)(W^T x_i - W^T \mu_i)^T$$

$$S_i^2 = W^T S_i' W$$

Combining all classes, total Variance =

$$S_1^2 + S_2^2 + S_3^2 + \cdots = W^T (S_1' + S_2' + \cdots) W$$

$$= W^T S_W W$$
$\uparrow$
Matrix

where

$$S_W = \sum_i S_i' = \sum_{\substack{\text{all data points} \\ x}} (x_i - \mu_i)(x_i - \mu_i)^T$$
$\qquad\qquad$ data $\quad$ mean of class it
$\qquad\qquad$ point $\qquad$ belongs to

Our Objective becomes

$$J(W) = \frac{W^T S_b W}{W^T S_W W} \quad \left( \text{form of generalized Rayleigh quotient} \right)$$

Differentiating & setting to 0 gives

$$S_b W = \lambda S_W W$$

$\left[ \begin{array}{l} S_b, S_W \text{ are known as} \\ \text{they can be calculated} \\ \text{from dataset} \end{array} \right]$

$\uparrow$
eigen values to be calculated

From this, W is found out

─── Note ───

Suppose there are two classes A & B
of people with height $x$ & weight $Y$

In general $\quad \underset{A}{Cov(X, Y)} \neq \underset{B}{Cov(X, Y)}$

But LDA assumes that

$$\underset{A}{Cov(X, Y)} = \underset{B}{Cov(X, Y)} = \underset{\text{Populations}}{Cov(X, Y)}$$

Finding the within class scatter takes
$N \times (d + d^2)$ time

Eigen decomposition & matrix multiplication
takes $O(d^3)$

If $N > d$ then     $O(d^3)$
   else            $O(N d^2)$ (if no of features
                              is trivial as compared to
                              samples)

Assumptions in LDA —

LDA assumes that data is normally
distributed within each class.

LDA assumes that every class has equal
covariance matrix

LDA is linear

Disadvantages —

① LDA is sensative to outliers
② LDA requires large number of samples
   relative to number of features


Classification →

The line obtained from LDA can not only be used
for data reduction, but also for classification
(eg using gaussian assumption or simple nearest neighbour)

After LDA projects the data, classification
can be done by various methods


eg→ Gaussian distribution assumption

Assuming that every distribution has a
gaussian and then like GMM finding the
class with maximum probability

$$\delta_k(x) = x^T \varepsilon^{-1} m_k - \frac{1}{2} m_k^T \varepsilon^{-1} m_k + \log P_{i_k}$$

(linear score function depends linearly on $x$)