

Week-11LDA (Linear Discriminant Analysis)Multi-class & binary class

2-class

$$\Sigma_W^{-1} \Sigma_B W = \lambda W$$

↳ Rank 1 Matrix

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

There is only one Eigen vector

$$\begin{matrix} \Sigma_W^{-1} & (\mu_1 - \mu_2) & (\mu_1 - \mu_2)^T & W \\ N \times N & N \times 1 & 1 \times N & N \times 1 \end{matrix}$$

↓

Let say Scalar Value = C

$$\Sigma_W^{-1} (\mu_1 - \mu_2) C = \lambda W$$

$$\Sigma_W^{-1} (\mu_1 - \mu_2) = \frac{\lambda}{C} W$$

$$\Sigma_W^{-1} (\mu_1 - \mu_2) \propto W$$

No matter what is N, there is only eigen vector for 2 class separation.

N = dimension of feature vector.



## LDA for k-classes

- In 2 classes, projection is done onto vector  $w$
- In k-class projection is done onto vector  $W$

$$W = [w_1, w_2, w_3 \dots w_{k-1}] \quad \text{where } w_i \text{ is } N \times 1$$

→ Projected means  $w^T \mu_1, w^T \mu_2 \dots w^T \mu_k$ .

→ Projected cov  $w^T \Sigma_1 w, w^T \Sigma_2 w \dots$

$$\begin{aligned} W^T \Sigma_1 W &= w_1^T \Sigma_1 w_1 + w_2^T \Sigma_1 w_2 \dots + w_{k-1}^T \Sigma_1 w_{k-1} \\ &= \text{Tr}(W^T \Sigma_1 W) \end{aligned}$$

$$\begin{aligned} \text{Similarly } w^T \Sigma_2 w &= \text{Tr}(w^T \Sigma_2 w) \\ W^T \Sigma_k W &= \text{Tr}(W^T \Sigma_k W) \end{aligned}$$

$$\text{Tr}(W^T \Sigma W) = \text{Tr}(W^T \Sigma_1 W) + \text{Tr}(W^T \Sigma_2 W) \dots \text{Tr}(W^T \Sigma_k W)$$

$\Sigma_w$  is within class covariance

Objective minimize  $\text{Tr}(W^T \Sigma_w W)$

$\Sigma_{B1} = (\mu_1 - \mu_2)(\mu_1^T - \mu_2^T)$  is b/w class cov.

$$E = W^T (\mu_1 - \mu_2)(\mu_1^T - \mu_2^T) W$$



$$E = \text{Tr}(W^T \Sigma_B W)$$

Goal is to maximize  $W^T \Sigma_B W$

$$\Sigma_B = \Sigma_{B1} + \Sigma_{B2} + \dots + \Sigma_{Bk}$$

$$\max \frac{\text{Tr}(W^T \Sigma_B W)}{\text{Tr}(W^T \Sigma_W W)}$$

$$L(W, \lambda) = (W^T \Sigma_B W) - \lambda (W^T \Sigma_W W - 1)$$

$$\frac{\partial L}{\partial W} = 2 \Sigma_B W - 2 \lambda \Sigma_W W = 0$$

$$\Sigma_B W = \lambda \Sigma_W W$$

$$\underbrace{\Sigma_W^{-1} \Sigma_B W}_{\downarrow} = \lambda W$$

~~W~~ is  $k-1$  rank matrix

thus have  $k-1$  eigen ~~vec~~ vectors

$$\Sigma_B = \begin{matrix} \Sigma_{B11} & \Sigma_{B12} & \dots & \Sigma_{B1k} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{Bk1} & \Sigma_{Bk2} & \dots & \Sigma_{B(k-1)k} \end{matrix}$$

With class

$$\Sigma_W = \Sigma_{B11} \oplus \Sigma_{B22} \oplus \dots \oplus \Sigma_{Bkk}$$



we know  $\Sigma_w$   
 we know  $\Sigma_T$  (Total variance)

$$\Sigma_T = \Sigma_w + \Sigma_B$$

$$\Sigma_B = \Sigma_T - \Sigma_w$$

$$\Sigma_T = \frac{1}{M} \sum_{i=1}^M (x_i - \mu)(x_i - \mu)^T$$

Pair-wise classifier

~~1 vs rest type~~

Thus. for  $k$ -classes, then  $k(k-1)/2$   
 pairwise classifiers

For test eg.

- Classification result from each classifier
- Final decision based on best classification among all the classifiers

LDA One vs rest

- $k$  classifiers in total
- Final decision based on best classifier among all the classifiers



## Bayesian Approach

Probability : Frequency count v/s Bayesian Approach

Conditional Probability

$P(A/B)$  = Prob. of happening A when B happened

$$P(X/Y) P(Y) = P(Y/X) P(X)$$

$$P(X/Y) = \frac{P(Y/X) \cdot P(X)}{P(Y)}$$

Core of Bayesian Approach.

Objective freq count  $\rightarrow$  # favourable / # Total

Subjective Bayesian: compute prob. in terms of other prob. values

Bayesian learning : for approximation

Eg. Regression

Bayesian Classification

- Bayes Optimal Classifier
- Naive Bayes
- Gibbs Sampling



## Bayesian Network

→ tree-like classifier using Bayesian Approach

$$P(Y/X) = P(X/Y)$$

## Bayesian Classification

$$P(Y/X) = \frac{P(X/Y) P(Y)}{P(X)}$$

X:- represent feature vector

Y:- Class or label

→  $P(Y/X)$  is a posteriori probability

Given X probability of represent Y

→  $P(X/Y)$  is likelihood prob or Class Condition, how well model represents a given set of feature vector X.

for given class label how what is prob X represent Y

→ prior prob:- a priori knowledge about model parameter Y  
 $P(Y)$

→ An evidence prob evidence about X  
 $P(X)$



Discriminative      Generative

$P(Y|x)$       Training      Testing

$P(x|Y)$       Testing      Training

for Generative Classifiers

Training       $P(x|y)$   
Testing       $P(y|x)$

$$P(x) = \sum y_i P(x|y_i) P(y_i)$$

Ideally

$$P(y) = P(x) = 1/N \quad N = \text{num class}$$

$$P(x|y) = P(y|x)$$

Same is exploited in Bayesian Approach

Freq count employs maximum likelihood approach for model parameter estimation

$$\hat{y} = \arg \max_y P(x|y)$$

Bayesian method employs posterior Approach for model parameter estimation

$$\hat{y} = \arg \max_y P(y|x)$$



## Bayesian Learning

Learning hypothesis function by exploiting bayesian concept in Bayesian learning

$y \rightarrow$  replaced by  $h$   
 $x \rightarrow D$

$$p(h/D) = \frac{p(D|h) \cdot p(h)}{p(D)}$$

How to find  $h$ ? (in Bayesian learning framework)

$$h_{MAP} = \max_{h \in H} p(h/D)$$

Using Bayes theorem  $h_{MAP} = \max_{h \in H} \frac{p(D|h) p(h)}{p(D)}$

$p(D)$  is independent of hypothesis

MAP: Maximum a posteriori based hypothesis

$$h_{MAP} = p(D|h) \cdot p(h)$$

If all hypothesis are equally likely i.e. prior probs  $p(h)$  are equal

we get  $h_{ML} = \max_{h \in H} p(D|h)$

This is called likelihood hypothesis  
also called Maximum likelihood (ML)



Finding hypothesis by Bayesian learning

Given function  $d_i = f(x_i) + \epsilon_i$

$d_i$  is from Normal distribution  $N(f(x_i), \sigma^2)$

let  $h$  represent hypothesis function represent  $f$

Can be estimated using Bayesian learning

$$h_{ML} = \arg\max_h p(D/W)$$

$$= \arg\max_h \prod_{i=1}^M \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{d_i - h(x_i)}{\sigma}\right)^2}$$

$$= \log\left(\arg\max_h \prod_{i=1}^M \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{d_i - h(x_i)}{\sigma}\right)^2}\right)$$

$$= \arg\min_h \sum_{i=1}^M \frac{1}{2} \left(\frac{d_i - h(x_i)}{\sigma}\right)^2$$

Linear least square error equation (Not prob. Approach)

Assumption: we have been of Generative process we have hypothesis function

Sample space of hypothesis is given.

hypothesis by non-probabilistic Approach is same as what we got in Bayes learning.



## Bayes Optimal Classifier

MAP may not be Optimal Classifier for a given test data

binary classification task (+ or -)

for eg. There are 3 hypothesis  $h_1, h_2, h_3$

MAP  $h_1$  has  $p = 0.4$  +  
 $h_2, h_3$  has  $p = 0.3$  - total  $0.6 > 0.4$

Then Bayes Optimal Classifier is useful

Given by  $\text{argmax}_j = \sum_{v_i \in V} p(v_j/h_i) p(h_i/D)$

→ It is called Optimal bcoz, no other classifier with same set of hypothesis and prior can outperform this on average.

→ hypothesis space is big & can't be applied in exhaustive way.

→ Gibbs Sampling that randomly chooses an hypothesis and finds class Error by this. It is within allowable limits.