

U 8941-5490

Ans to the Que No 1(a)

Quiz-4

Date :

Ans to Que 1(a)

Posterior-ratio :

$$\frac{P(c_n=1|x_n)}{P(c_n=2|x_n)} \Rightarrow \frac{P(x_n|c_n=1)P(c_1)}{P(x_n|c_n=2)P(c_2)} \begin{matrix} > 1 \\ < 1 \end{matrix}$$

if the posterior-ratio is > 1 then c_1 if the ratio is < 1 then c_2

$$\frac{P(x_n|c_n=1)}{P(x_n|c_n=2)} \begin{matrix} > \\ < \end{matrix} \frac{P(c_2)}{P(c_1)} \cdot 1 \Rightarrow \text{threshold}$$

$$\Rightarrow \log P(x_n|c_n=1) - \log P(x_n|c_n=2) \begin{matrix} > \\ < \end{matrix} \log \frac{P(c_2)}{P(c_1)}$$

$$\Rightarrow \log \mathcal{N}(x_n | \mu_1, \Sigma_1) - \log \mathcal{N}(x_n | \mu_2, \Sigma_2) \begin{matrix} > \\ < \end{matrix} \log \frac{P(c_2)}{P(c_1)}$$

(i)

(2)

Date :

$$\log N(x|\mu, \Sigma) = \log \left[\frac{\exp\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\}}{(2\pi)^{d/2} |\Sigma|^{1/2}} \right]$$

$$= -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma|$$

(1)

log of a general multi-variate
Gaussian PDF

so, we can re-write eqⁿ (i) as

$$-\frac{1}{2} (x_n - \mu_1)^T \Sigma_1^{-1} (x_n - \mu_1) - \frac{d}{2} \log 2\pi$$

$$+ \frac{1}{2} \log |\Sigma_1^{-1}| + \frac{1}{2} (x_n - \mu_2)^T \Sigma_2^{-1} (x_n - \mu_2)$$

$$+ \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_2^{-1}|$$

From here on we will just work
with the ~~L.H.S~~ of the above eqⁿ.

(3)



Date :

Again, $-\frac{1}{2} (X - \mu)^T \Sigma^{-1} (X - \mu)$

$$= -\frac{1}{2} \left[X^T \Sigma^{-1} X - X^T \Sigma^{-1} \mu - \mu^T \Sigma^{-1} X + \mu^T \Sigma^{-1} \mu \right]$$

Accordingly we can rewrite +

$$= -\frac{1}{2} \left[X^T \Sigma^{-1} X - 2 X^T \Sigma^{-1} \mu + \mu^T \Sigma^{-1} \mu \right]$$

So, we can rewrite the eqⁿ of the previous page as :

$$\frac{1}{2} \left\{ X_n^T (\Sigma_2^{-1} - \Sigma_1^{-1}) X_n - 2 X^T (\Sigma_2^{-1} \mu_2 - \Sigma_1^{-1} \mu_1) \right. \\ \left. + \text{const} \left\{ \sum_{c_2}^{c_1} \log t \right\} \right.$$

$$\Rightarrow \gamma(u) = X^T A^0 X + b^T X \sum_{c_2}^{c_1} 2 \log t$$

if $\gamma > t$ then c_1
 \circ/ω c_2 .



(4)

Date :

$$\text{So, } A = (\Sigma_2^{-1} - \Sigma_1^{-1})$$

$$b = (\Sigma_2^{-1} \mu_2 - \Sigma_1^{-1} \mu_1)$$

$$k = \frac{P(c_2)}{P(c_1)}$$

* we can ignore the constant as it is same for every data instance X_n .

$$X^T X = \lambda I$$

Ans to the Que No 1(b)

5

Date :

Ans to the Que 1(b)

LDA from QDA :

Assuming, $\Sigma_1 = \Sigma_2 = \Sigma$ that is the covariance matrix of the two classes are same. which means the ~~two~~ two curves have the same variance. The contours of the two gaussian curve on 2D space has the same shape.

$$\text{Then, } A = (\Sigma_2^{-1} - \Sigma_1^{-1}) (\Sigma_1^{-1} - \Sigma_2^{-1}) = 0$$

So the class QDA becomes,

$$y(x) = b^T X = \textcircled{X^T b}$$

(6)

Date :

where, $b = \Sigma^{-1}(\mu_2 - \mu_1) = w$

So, $-2x^T \Sigma^{-1} \cdot (\mu_2 - \mu_1) \sum_{c_2}^{c_1} 2 \log +$

$\Rightarrow -2x^T \Sigma^{-1} (\mu_2 - \mu_1) \sum_{c_2}^{c_1} 2 \log \frac{P(c_2)}{P(c_1)}$

$\Rightarrow x^T \Sigma^{-1} (\mu_2 - \mu_1) \sum_{c_2}^{c_1} - \log \frac{P(c_2)}{P(c_1)}$

$\Rightarrow x^T \Sigma^{-1} (\mu_2 - \mu_1) \sum_{c_2}^{c_1} \log \frac{P(c_1)}{P(c_2)}$

$\therefore h = \frac{P(c_1)}{P(c_2)}$

$b = \Sigma^{-1} (\mu_2 - \mu_1)$

Ans to the Que No 1(c)

(7)

Date :

Ans to the Que 1(c)

If it is assumed that the covariance matrix is diagonal then ^{from} QDA we can obtain Gaussian Naive Bayes classifier. That means, ~~Every~~ the features are independent of each other.

In that case we can assume a univariate model.

$$X_{nd} \sim \mathcal{N}(\mu_d, \sigma_d^2)$$

$$P(X_n | c_n) = \prod_{i=1}^d P(X_{ni} | c_n)$$

d = dimensionality of the features.

$X_{ni} \rightarrow$ Refers to the i 'th feature of the n 'th data instance.



Extra Credit:

(8)

Date:

$$\frac{P(c_n=1 | x_n)}{P(c_n=2 | x_n)} = \frac{\prod_{i=1}^d P(x_{ni} | c_n=1) P(c_n=1)}{\prod_{i=1}^d P(x_{ni} | c_n=2) P(c_n=2)}$$

$$\Rightarrow \frac{\log(c_n=1 | x_n)}{\log(c_n=2 | x_n)} = \frac{\log \prod_{i=1}^d P(x_{ni} | c_n=1) P(c_n=1)}{\log \prod_{i=1}^d P(x_{ni} | c_n=2) P(c_n=2)}$$

$$\Rightarrow \frac{\sum_{i=1}^d \log P(x_{ni} | c_n=1)}{\sum_{i=1}^d \log P(x_{ni} | c_n=2)} < \frac{\sum_{c_1} \log P(c_2)}{\sum_{c_2} \log P(c_1)} \Rightarrow +$$

$$\Rightarrow \sum_{i=1}^d \log N(x_{ni} | \mu_i, \sigma_i^2) - \sum_{i=1}^d \log N(x_{ni} | \mu_i, \sigma_i^2)$$

$$\frac{\sum_{c_1} \log +}{\sum_{c_2} \log +}$$

$$N(\mu_1, \sigma_1^2) \Rightarrow c_1 \quad N(\mu_2, \sigma_2^2) \Rightarrow c_2$$

9

Date :

$$\Rightarrow \sum_{i=1}^d -\log \mathcal{N}(x_{ni} | \mu_i, \sigma^2) - \sum_{i=1}^d \log \mathcal{N}(x_{ni} | \mu_i, \sigma^2)$$

$$N(X | \mu_1, \sigma^2) \Rightarrow C_1$$

$$N(x | \mu_2, \sigma^v) \Rightarrow c_2$$

$$\log \mathcal{N}(x | \mu, \sigma^2) = -\frac{1}{2\sigma^2} (x_n - \mu)^2 - \frac{1}{2} \ln \sigma^2 - \frac{1}{2} \ln 2\pi$$

So we can rewrite eqn. (i) as

$$= \sum_{i=1}^n \left[-\frac{1}{2\sigma_i^2} (x_{ni} - \mu_i)^2 - \frac{1}{2} \ln \sigma_i^2 - \frac{1}{2} \ln 2\pi \right]$$

$$\sum_{i=1}^n \log t$$

$$\Rightarrow \sum_{i=1}^d \left[\frac{1}{\sigma_i} \cancel{x_{ni}} - \frac{1}{\sigma_i} 2 x_{ni} \mu_i + \mu_i \frac{1}{\sigma_i} \cancel{x_{ni}} \right] - \frac{1}{\sigma_i} \cancel{x_{ni}} + 2 x_{ni} \mu_i \frac{1}{\sigma_i} - \frac{1}{\sigma_i} \mu_i^2$$

$$\sum_{i=1}^d \log t$$


$$\Rightarrow \sum_{i=1}^d \frac{1}{2 \sigma_i} \left[2 x_{ni} (\mu_i - \mu_i) \right] \log t$$

$$\Rightarrow \underbrace{\sum_{i=1}^d \frac{1}{\sigma_i} x_{ni} (\mu_i - \mu_i)}_{w^T \bar{x}} + \underbrace{\sum_{i=1}^d \frac{1}{2 \sigma_i} (\mu_i - \mu_i)}_{w_0} \log t$$

$$\Rightarrow \boxed{w^T \bar{x} + w_0 \sum_{i=1}^d \log t}$$

It is a linear classifier

Ans to the Que No 2



Date :

Ans to the Que No 2

Objective

- Maximize the margin between Decision boundary and data points.
- It is a decision machine so no posterior probability.
- $$y(x) = w^T \phi(x) + b$$

↓
non-linear function
- Main objective is to maximize the margin between support vectors and decision boundary.
- we want to find w and b that gives an maximum margin.

(12)

Date :

kernel trick :

$$\rightarrow \boxed{f(u) = w^T \phi(u) + b} \Rightarrow \text{it is a linear}$$

function. But in order to solve the non-linear problems we have $\phi(u)$.

which is a non-linear feature extraction function. That means using

ϕ we are transforming/mapping data from low ~~to~~ to high dimension. On

that higher dimensional feature space the ~~data become~~ problem becomes linear.

\rightarrow We don't need to know the exact

ϕ . Using kernel trick we can achieve this non-linearity.

\rightarrow Using the kernel trick SVM avoids solving/optimizing for learning the best ϕ function.

Training

→ During training SVM procedure SVM considers a particular ϕ choice.

→ So at the beginning SVM is restricted to a ϕ choice.

→ We want to minimize classification error. \Rightarrow Maximize # of correctly classified points.

$$\Rightarrow \max \frac{\sum_n y(x_n)}{\|w\|} = \max \frac{\sum_n (w^T \phi(x_n) + b)}{\|w\|}$$

$$y(x_n) = w^T \phi(x_n) + b$$

for correctly classified points

$$y(x_n) > 0$$

where, $y = \begin{cases} +1 \\ -1 \end{cases}$

~~class~~
 ~~label~~

14

Date :

→ So our target is to find the w and b that maximizes the ~~low~~ margin for the support vector. As our target is to find a unit solⁿ for w we have $\text{norm}(w) \Rightarrow \|w\|$ in the denominator.

$$\max \frac{\ln y(x_n)}{\|w\|} \quad \max \quad \text{minimum margin is } 1$$

→ So we now fix $\ln y(x_n) \geq 1$

and will use lagrange multiplier to optimize for w and b . So minimize the denominator which is $\|w\|$.

→ using lagrange multiplier we convert this constrained problem to an unconstrained problem.

15

Lagrange mult.

Date :

$$w = \sum_{n=1}^N \alpha_n \phi(x_n)$$

* now we search for α_n .

$$\rightarrow f(x_n) = \sum_{n=1}^N w^T \phi(x_n) + b$$

$$= \sum_{n=1}^N \alpha_n \phi(x_n)^T \phi(x_n) + b$$

Kernel function

kernel trick

 $N \Rightarrow$ # number of training samples.

→ In the beginning we ~~not~~ started with the S.V only. Later we see that we are including all the training samples (including S.V.) while considering the loss function later we will see that eventually only the S.V ~~not~~ has the contributes to maximizing the margin. And also we don't always know what are the S.V beforehand.

Testing:

$$\rightarrow f(x) = \sum_{n=1}^N \alpha_n t_n k(x, x_n)$$

test point

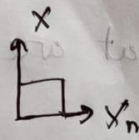
training point

→ So to decide for a test point x we take the inner-product (similarity) between x and all the training points x_n .

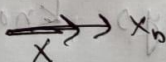
→ So the similarity score should be high if x and x_n is similar. So

the kernel function works as a weighting factor for the ^{train data} label


~~if they are similar~~



$$x_n^T x = 0$$



$$x_n^T x = 1 \rightarrow \text{max}$$



$$x_n^T x = -1 \rightarrow \text{min}$$

17

Date :

→ If x and x_n are similar we count that label t_n in decision making

→ if $\boxed{x_n^T x \rightarrow 0}$ then t_n has no contribution

→ if $\boxed{x_n^T x = -1}$ then t_n contributes negatively in decision making.

Loss function:

$$\rightarrow E(y_n, t_n) = \max\{1 - y_n t_n, 0\}$$

$y_n t_n \gg 1$ for every support training point other than the support vectors so, the error for them is 0.

But the error for support vectors would be non-zero as $y_n t_n < 1$

18

Date :

for S.V. That means we want to
 maximize the margin for the S.V.
 The margin is already maximized
 for other training samples.

$$K = X^T X$$

Ans to the Que No 3

Boosting:

- Instead of using only one model. Data is trained on multiple models.
- Each new model is trained to emphasize the data that was incorrectly classified by the previous model.
- Uses some loss functions to measure the performance of a model.
- Size of trees is a parameter that needs to be tuned.

Bagging:

- Each model has equal weight in voting.
- But each model is trained on a randomly sampled subset of the original train data. This is done to avoid overfitting to train data and to promote model variance. Example: Random Forest => Feature Bagging as RF selects a random subset of features within the random subset of the dataset.
- Number of trees to be built is a parameter that needs tuning.
Normally uses out of bag error or any loss function to measure the optimal value of number of trees to be built.
- Requires mechanism to de-correlate the trees that are built => bootstrap sampling is one such mechanism.