

# Lecture 5<sup>LDA</sup>, Bayesian decision theory

Note Title

9/26/2017

$X \xrightarrow{f} y$  estimate  $f$  or estimate  $P(y|x)$  or  $P(x,y)$

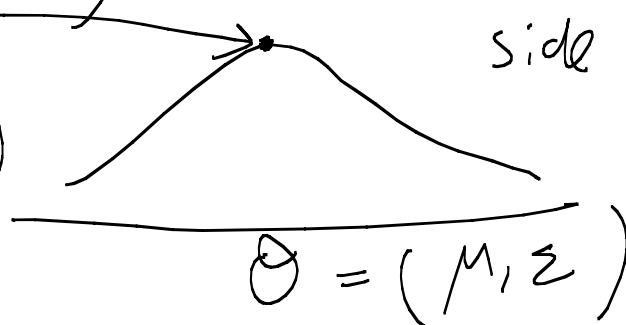
$$P(y=c|x,\theta) = \frac{P(x|y=c,\theta) P(y=c|\theta)}{\sum_c P(x|y=c,\theta) P(y=c|\theta)}$$

$$D = \{ (x_i, y_i) \} \sim N(\mu, \Sigma)$$

MLE, MAP  $\leftarrow P(\theta|D)$  in Bayesian setting

$$\frac{P(D|\theta) \text{ or } \log P(D|\theta)}$$

$$\frac{P(D|\theta) P(\theta)}{N}$$



$$P(y=c|x,\theta) = \frac{P(x|y=c,\theta) P(y=c|\theta)}{\sum_c P(x|y=c,\theta) P(y=c|\theta)}$$

classifier

$$\sum_c P(x|y=c,\theta) P(y=c|\theta) \xrightarrow{\text{joint}}$$

generative classifier

↓  
discriminative classifier will have a  
direct function

Back to modelling class conditional densities.

$$\hookrightarrow P(X|Y=c)$$

Naive  
Bayes classifier

$$D. \prod_{i=1}^D P(x_i|Y=c)$$

$$\cancel{P}(X|Y=c) =$$

$$\sim N(\mu_c, \Sigma_c)$$

= Gaussian

$\Sigma = \Sigma_c$  discriminant  
analysis

otherwise

QDA  
Quadratic discriminant

②

LDA linear  
discriminant analysis

|

Back to naive Bayes

(2)

$$\Sigma = \Sigma_c = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_d^2 \end{bmatrix}$$

$$= \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}}$$

look into

$$\exp\left(-\frac{1}{2} (X-\mu)^T \Sigma (X-\mu)\right)$$

$$= \frac{1}{(2\pi)^{D/2} \prod_{i=1}^D (\sigma_i^2)^{1/2}} \exp\left(-\frac{1}{2} \left( \sigma_1^2 (x_1 - \mu_1)^2 + \sigma_2^2 (x_2 - \mu_2)^2 + \dots + \sigma_d^2 (x_d - \mu_d)^2 \right)\right)$$

$$= \exp\left(-\frac{1}{2} \sigma_1^2 (x_1 - \mu_1)^2\right) \exp\left(-\frac{1}{2} \sigma_2^2 (x_2 - \mu_2)^2\right) \dots \exp\left(-\frac{1}{2} \sigma_d^2 (x_d - \mu_d)^2\right)$$

Summary =  $N(\mu_1, \sigma_1) \dots N(\mu_D, \sigma_D)$

look into this to find decision boundary

$$\rightarrow P(Y=1|X) = P(Y=2|X)$$

$$\frac{1}{(2\pi)^{D/2}} \exp\left(-\frac{1}{2} (X-\mu_1)^T \Sigma_1 (X-\mu_1)\right) = \frac{1}{(2\pi)^{D/2}} \exp\left(-\frac{1}{2} (X-\mu_2)^T \Sigma_2 (X-\mu_2)\right)$$

4.2.9 MLE for discriminant analysis

$$\log P(D|\theta) \quad (\text{IID})$$


$$= \sum_{i=1}^N \sum_{c=1}^C \mathbb{I}(y_i=c) \log \pi_c + \sum_{c=1}^C \left[ \sum_{i: y_i=c} \log N(x_i | \mu_c, \Sigma_c) \right]$$

$$\hat{\pi}_c = N_c / N$$

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i: y_i=c} x_i = \frac{1}{N_c} \sum_{i=1}^N \mathbb{I}(y_i=c) x_i$$

$$\hat{\Sigma}_c = \frac{1}{N_c} \sum_{i: y_i=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$$

---



side  $\Sigma_c^\wedge$  need to be invertible  
i.e. it need to be full rank.  
 $= D$

$N_c < D \Rightarrow$  Can't invert

$N_c \gg D \Rightarrow$  high chances  
of being  
invertible

---

Even in ImageNet 2011 0.6 million  
big data exa when  
we cluster based on class  
label we ~~are~~ have  
less sample  
1000 classes  
So 1000 images / class

each image

$$X_{300 \times 300 \times 3} \in \mathbb{R}^{270000=1}$$

$$N_c = 1000 \ll D (270000)$$

What to do?

- Duplicate sample with perturbation
- Dimensionality reduction  
(PCA)
- use diagonal matrix ( $D$  parameter per class)
- $\Sigma_c = \Sigma$   $\forall c$   
or  $\Sigma_c = \Sigma$  and diagonal

— Use Full covariance but impose

side — determinant of matrix = product of eigenvalues prior  
— For a matrix sum of eigenvalue = sum of diagonal entries  
= trace of matrix

## Chapter 5

### 5.1 Bayesian Decision theory

$$X \sim N(\mu, \Sigma)$$

choose some action  $a \in A$  (action space)

How well we did?

define a loss  $L(y_{\text{true}}, a)$

and we want to minimize the loss

$$L(y_{\text{true}}, a) = \mathbb{I}(y_{\text{true}} \neq a) \quad [\text{mis classification loss}]$$

What about regression?

$$L(y_{\text{true}}, a) = (y_{\text{true}} - a)^2$$

Decision rule  $S: X \rightarrow A$

optimal policy

$$S(x) = \arg \min_{a \in A} \mathbb{E}[L(y, a)]$$

In Bayesian approach to decision theory, the optimal action, have  $x \in X$ , is defined as the action that minimizes the posterior Expected loss

$$p(a|x) = \mathbb{E}_{(y|x)}[L(y, a)] = \sum_y L(y, a) P(y|x)$$



Bayes estimate (Bayes decision rule)

$$\delta(x) = \arg \min_{a \in \mathcal{A}} p(a|x)$$

MAP estimate minimize 0-1 loss

$$L(y, a) = L(y, \hat{y}) = \mathbb{I}(y \neq \hat{y}) = \begin{cases} 0 & \text{if } y = \hat{y} \\ 1 & \text{if } y \neq \hat{y} \end{cases}$$

	$\hat{y} = 1$	$\hat{y} = 0$
$y = 1$	0	1
$y = 0$	1	0

Then posterior expected loss is

$$P(a|x) = 1 - P(a \neq y|x) = 1 - P(a \neq \hat{y}|x)$$

$$a = \hat{y} = \arg \max_{a \in \mathcal{Y}} p(a=y|x)$$

justifies our earlier choice of  
MAP

5.7.1.3

posterior mean minimizes  $L_2$   
loss for continuous parameter

$p=2$

$$\hat{y} = a = E[a|X] \quad [a-y]^p$$

If  $p=1$  basically  $|a-y|$

then optimal action is posterior

median 
$$P(y < a|X) = P(y > a|X) = .5$$

---

Note outlier become more prominent to the power  $p$ . like in image normalization using 11411p or max value may not be the right thing to do. Normalize using 98 or 99 percentile.

---

5.7.2

The false positive and

false negative

Binary setting

	$\hat{y}=1$	$\hat{y}=0$
$y=1$	0	$L_{FN}$
$y=0$	$L_{FP}$	0

Posterior Expected loss

$$P(\hat{y}=0|X) = L_{FN} P(y=1|X)$$

$$P(\hat{y}=1|X) = L_{FP} P(y=0|X)$$

Hence pick  $\hat{y}=1$  if

$$P(\hat{y}=0|X) > P(\hat{y}=1|X)$$

$$L_{FN} P(y=1|X) > L_{FP} P(y=0|X)$$

$$\frac{P(y=1|X)}{P(y=0|X)} > \frac{L_{FP}}{L_{FN}} = \tau$$

$$\underline{L_{FP} = C L_{FN}}$$

5.7.2.1      ROC curves , AUC

Suppose we are in supervised  
binary decision setting  $D = \{(x_i, y_i)\}$   
we apply  $\Gamma$  and count no of  
true, positive, false, positive, true  
negative, False negative

		True		Total
		1	0	
Estimate	1	TP	FP	$\hat{N}_+ = TP + FP$
	0	FN	TN	$\hat{N}_- = TN + FN$

$$\text{Total} \quad \left[ N_+ = TP + FN = \frac{N_-}{FP + TN} \right]$$

$$TPR (\text{sensitivity, recall, hit rate}) = \frac{TP}{N_+}$$

$$FPR (\text{False alarm, type I error}) = \frac{FP}{N_-}$$

ROC plot of TPR vs FPR for various value of  $\tau$

(Receiver operating characteristic)

$$\tau = 1 \Rightarrow FPR = 0 \neq TPR = 0$$

$$\tau = 0 \Rightarrow FPR = FPR = 1$$

ROC is summarized using one number called AUC (Area under the curve)

