

## 2.13 Mutual Information $I(X_1, X_2)$

(3 point)  $I(X_1, X_2) = H(X_1) + H(X_2) - H(X_1, X_2)$

where  $H$  is entropy of a random variable.

Hence as per given values for entropies in book

$$= \frac{1}{2} \log_2 [2\pi e \sigma^2] + \frac{1}{2} \log_2 [2\pi e \sigma^2] - \frac{1}{2} \log_2 [(2\pi e)^2 (\sigma^4 - \sigma^4 \rho^2)]$$

$$= \log_2 [2\pi e \sigma^2] - \frac{1}{2} \log_2 [(2\pi e)^2 \sigma^4 (1 - \rho^2)]$$

$$= \log_2 \frac{2\pi e \sigma^2}{[(2\pi e)^2 \sigma^4 (1 - \rho^2)]^{1/2}}$$

$$= \log_2 \frac{\cancel{2\pi e} \sigma^2}{\cancel{2\pi e} \sigma^2 (1 - \rho^2)^{1/2}}$$

$$= \log_2 \frac{1}{(1 - \rho^2)^{1/2}}$$

(take coefficient of second log back into power and using  $\log \frac{a}{b} = \log a - \log b$ )

If  $\rho = \pm 1$  or  $-1$  ( $X_1, X_2$  highly correlated)

then  $I(X_1, X_2) = \infty$  justifies this

If  $\rho = 0$  i.e.  $X_1, X_2$  are statistically not correlated at all then

$I(X_1, X_2) = \log 1 = 0$  justifying this.

①

### 3.6 MLE For the Poisson distribution

(Points 4)

Poisson pmf  $Poi(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$  [note the support in 0, 1, 2, ... all natural number]

Hence let  $D = \{x_i\}_{i=1}^N$  be our  $i.i.d$  observations from this distribution

then log likelihood of  $D$  is

$$= \log P(D) = \log \prod_{i=1}^N Poi(x_i|\lambda) = \sum_{i=1}^N \log Poi(x_i|\lambda)$$

$$= \sum_{i=1}^N \log e^{-\lambda} \frac{\lambda^{x_i}}{x_i!}$$

$$= \sum_{i=1}^N (-\lambda) + \sum_{i=1}^N x_i \log \lambda - \sum_{i=1}^N \log(x_i!)$$

$$= -N\lambda + \sum_{i=1}^N x_i \log \lambda - \sum_{i=1}^N \log(x_i!)$$

Taking derivative with respect to  $\lambda$  and equating to zero yield

$$-N + \frac{\sum_{i=1}^N x_i}{\lambda} = 0 \Rightarrow \lambda = \frac{\sum_{i=1}^N x_i}{N}$$

### 3.7

$$P(\lambda|D) \propto P(D|\lambda) P(\lambda)$$

(a)  
(2 points)  
each  
(a, b, ...)

$$\propto \left( \prod_{i=1}^N Poi(x_i|\lambda) \right) Gra(\lambda|a, b)$$

(using given conjugate prior  $Gra(\lambda|a, b)$  for  $\lambda$ )

$$\propto \left( \prod_{i=1}^N e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \right) \left( \lambda^{a-1} e^{-\lambda b} \right)$$

$$\propto \frac{\lambda^{\sum_{i=1}^N x_i + a - 1} e^{-N\lambda - b\lambda}}{\prod_{i=1}^N x_i!}$$

$$= Gra\left(\lambda \mid a + \sum_{i=1}^N x_i, b + N\right)$$

(b) Posterior mean is  $E[\lambda|D] = \frac{a + \sum_{i=1}^N x_i}{b + N}$

as  $a \rightarrow 0, b \rightarrow 0$ , it tends to MLE estimate  $\frac{\sum_{i=1}^N x_i}{N}$  (without prior on  $\lambda$ )

3.11 (only for ENCE 4630) For log likelihood proceed as in previous problem, ②  
 (Point each)  
 (a) it should be

$$L(\theta) = N \log \theta - \theta \sum_{i=1}^N x_i$$

For MLE estimate of  $\theta$ , take derivative and equate to zero

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{N}{\theta} - \sum_{i=1}^N x_i \Rightarrow \hat{\theta} = \frac{N}{\sum_{i=1}^N x_i}$$

(b)  $\hat{\theta}_{MLE}(D) = \hat{\theta}_{MLE}(\{x_1=5, x_2=6, x_3=4\})$  ( $\because D = \{5, 6, 4\}$ )

$$= \frac{3}{5+6+4} = \frac{3}{15} = \frac{1}{5}$$

(c)  $E(\theta) = 1/\lambda = \frac{1}{3} \Rightarrow \hat{\lambda} = 3$  ( $\because p(\theta) \propto \theta^{-1} e^{-\theta \lambda}$ )

$-\theta \lambda = \text{Expon}(\theta | \lambda)$   
 $= \text{Gamma}(\theta | 1, \lambda)$   
 Hence exponential is a special case of gamma

(d) posterior  $p(\theta | D, \hat{\lambda}) \propto p(D | \theta) p(\theta | \hat{\lambda})$

$$\propto \theta e^{-\theta 5} \theta e^{-\theta 6} \theta e^{-\theta 4} e^{-3\theta}$$

$$\propto \theta^3 e^{-18\theta} \propto \theta^{4-1} e^{-18\theta} = \text{Gamma}(\theta | 4, 18)$$

( $\because \text{Gamma}(\theta | a, b) \propto \theta^{a-1} e^{-b\theta}$ )

(e) yes, exponential is a special case of Gamma.

(f) posterior  $p(\theta | D, \hat{\lambda})$  is  $\text{Gamma}(\theta | 4, 18)$   
 Hence mean =  $4/18$

(g) posterior mean ( $4/18$ ) is a probabilistic adjustment between prior mean ( $1/3$ , expert choice) and data driven mean ( $\text{MLE} = 1/5$ ). In small sample size we should take expert advice. Hence posterior

mean is reasonable.

(3)

3.20 (1 point each)

(a) since features are not conditionally independent we need to specify probability for each configuration of bit in  $X \in \{0,1\}^D$

Hence we need  $2^D - 1$  parameters/class

or  $C(2^D - 1)$  parameter for classes

we would need  $C$  different probability distributions  $P(X|C)$  or histogram. [One choice can be

binomial  $(K; D, \theta)$ ]

[probability of  $K$  success in  $D$  trial of coin with parameter  $\theta$ ]

Note: We cannot use multivariate Gaussian as data is binary vector.

(b) Naive Bayes based model will work better as it has less number of parameters. With less sample size it will not overfit and parameters estimation will be more reliable.

(c) with large sample size, use full model as it is more accurate and there is enough data to reliably estimate parameters.

(d) Both take  $O(ND)$  time. As estimation (4) of parameters is mostly about counts and we can do this by iterating over  $N$  examples and updating counts.

(e) Both model take  $O(D)$  time

$$\text{For NB} \quad P(y=c|x, \theta) = \prod_{j=1}^D \theta_{j,c}^{\mathbb{I}(x_j=1)}$$

We know  $\theta_{j,c}$  after estimation and indexing  $x_{\text{test}} \in \{0,1\}^D$  ~~can be done in~~ to set  $x_{\text{test}}[j] = x_j$  can be done in  $O(D)$  time.

For Full model we need to pick the parameter corresponding to bit pattern in  $x_{\text{test}} \in \{0,1\}^D$ . This is equivalent to converting  $x_{\text{test}}$  to an integer and can be done in  $O(D)$  time. If we use hash it can be done in  $O(1)$  constant time.

(f) we can play the marginalizing trick from probability we need to compute

$$P(y|x_v, \hat{\theta}) \propto P(x_v|y, \hat{\theta}) P(y) \propto \sum_{x_n} P(x_v, x_n|y, \hat{\theta})$$

$$\begin{aligned} \text{for Naive Bayes} \quad \sum_{x_n} P(x_v, x_n|y, \hat{\theta}) &= \sum_{x_n} P(x_v|y, \hat{\theta}) P(x_n|y, \hat{\theta}) \\ &= P(x_v|y, \hat{\theta}) \left( \sum_{x_n} P(x_n|y, \hat{\theta}) \right) \end{aligned}$$

If we have pre-computed  $\sum_{x_n} (x_n | y, \hat{\theta})$  (5)

then in naive Bayes we need  $O(v)$  time

For Full model we need to enumerate over  
all  $2^h$  values of  $x_n \in \{0, 1\}^h$

for marginalising  $\sum_{x_n} P(x_v, x_n | y, \hat{\theta})$ .

Hence it takes  $O(D 2^h)$   $\approx O(2^h)$   
computational time if  $D \ll 2^h$