

# CS391L HW3: Problem Set

Mingyo Seo,  
UT EID: ms84662  
Email: mingyo@utexas.edu

## I. BASIC PROBABILITY

The **answer** is  $p(X|Y) = \frac{1}{73}$ . For the notation, please see the below description.

Let  $X$  be the case that it will actually rain today, and  $Y$  be the case that the meteorologist predicts a rainy day. In other words,  $X^c$  means that it will not rain today. In the case of  $Y^c$ , the meteorologist predicts that it would not rain today. From the given historical data,

$$p(X) = \frac{5}{365} = \frac{1}{73}, \quad (1)$$

$$p(X^c) = 1 - P(X) = \frac{72}{73}. \quad (2)$$

The meteorologist correctly predicts 90% when it rains, so

$$p(Y|X) = \frac{9}{10}, \quad (3)$$

The meteorologist correctly predicts 10% when it doesn't rain, so

$$p(Y^c|X^c) = \frac{1}{10}, \quad (4)$$

This also implies that

$$p(Y|X^c) = 1 - p(Y^c|X^c) = \frac{9}{10}. \quad (5)$$

Given the condition that a meteorologist predicts rain today, the conditional probability of rain today can be formulated as,

$$p(X|Y) = \frac{p(X, Y)}{p(Y)} = \frac{p(X, Y)}{p(X, Y) + p(X^c, Y)}. \quad (6)$$

From Bayes' theorem, we can get

$$p(X, Y) = p(Y|X)p(X), \quad (7)$$

$$p(X^c, Y) = p(Y|X^c)p(X^c). \quad (8)$$

Therefore, from Equation 1, 2, 3, 5, we can compute the probabilities,

$$p(X, Y) = \frac{9}{10} \times \frac{1}{73}, \quad (9)$$

$$p(X^c, Y) = \frac{9}{10} \times \frac{72}{73}. \quad (10)$$

Therefore, we can compute the **answer** by plugging these probabilities in 6, as

$$p(X|Y) = \frac{p(X, Y)}{p(X, Y) + p(X^c, Y)} = \frac{\frac{9}{10} \times \frac{1}{73}}{\left(\frac{9}{10} \times \frac{1}{73}\right) + \left(\frac{9}{10} \times \frac{72}{73}\right)} = \frac{1}{73}. \quad (11)$$

## II. ENTROPY

From the definitions,

$$H[\mathbf{y}] = - \int p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{y}, \quad (12)$$

$$H[\mathbf{y}|\mathbf{x}] = - \int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{y} d\mathbf{x}. \quad (13)$$

From Bayes' theorem, we can get

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}, \quad (14)$$

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y}). \quad (15)$$

Therefore, we can re-write the given KL definition of mutual informaiton, as

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &= - \int \int p(\mathbf{x}, \mathbf{y}) \ln \left\{ \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right\} d\mathbf{x} d\mathbf{y} \\ &= - \int \int p(\mathbf{x}, \mathbf{y}) \ln \frac{p(\mathbf{y})}{p(\mathbf{y}|\mathbf{x})} d\mathbf{x} d\mathbf{y} \\ &= - \int \int p(\mathbf{x}, \mathbf{y}) \{ \ln p(\mathbf{y}) - \ln p(\mathbf{y}|\mathbf{x}) \} d\mathbf{x} d\mathbf{y} \\ &= - \int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}) d\mathbf{x} d\mathbf{y} + \int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{x} d\mathbf{y}, \\ &= - \int \int p(\mathbf{y}) p(\mathbf{x}|\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{x} d\mathbf{y} + \int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{x} d\mathbf{y}. \end{aligned} \quad (16)$$

By using the property,

$$\forall \mathbf{y}, \int p(\mathbf{x}|\mathbf{y}) d\mathbf{x} = 1, \quad (17)$$

we can re-write the first term of the above equation with respect to  $H[\mathbf{x}]$  at Equation 12, as

$$\begin{aligned} - \int \int p(\mathbf{x}|\mathbf{y}) p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{x} d\mathbf{y} &= - \int p(\mathbf{y}) \ln p(\mathbf{y}) \left\{ \int p(\mathbf{x}|\mathbf{y}) d\mathbf{x} \right\} d\mathbf{y} \\ &= - \int p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{y} \\ &= H[\mathbf{y}]. \end{aligned} \quad (18)$$

The second term of Equation 16 can be re-written with respect to  $H[\mathbf{y}|\mathbf{x}]$  at Equation 13 as

$$\begin{aligned} \int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{x} d\mathbf{y} &= - \left( - \int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{x} d\mathbf{y} \right) \\ &= - \left( - \int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{y} d\mathbf{x} \right) \\ &= -H[\mathbf{y}|\mathbf{x}]. \end{aligned} \quad (19)$$

Therefore, we get the results,

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]. \quad (20)$$

Due to the symmetric form of  $I[\mathbf{x}, \mathbf{y}]$ , by switching  $\mathbf{x}$  and  $\mathbf{y}$  at Equation 16, we can also get

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &= I[\mathbf{y}, \mathbf{x}] \\ &= H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}]. \end{aligned} \quad (21)$$

### III. BETA DISTRIBUTION

From the definitions Beta distribution,

$$p(\mu) = \frac{\mu^{a-1} (1-\mu)^{b-1}}{C(a,b)}, \quad (22)$$

where  $C(a,b)$  is the normalizing constant, given as

$$\begin{aligned} C(a,b) &= \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu, \\ &= \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}. \end{aligned} \quad (23)$$

Then,  $E(\mu)$  is given as,

$$\begin{aligned} E(\mu) &= \int_0^1 \mu p(\mu) d\mu \\ &= \int_0^1 \mu \frac{\mu^{a-1} (1-\mu)^{b-1}}{C(a,b)} d\mu \\ &= \frac{1}{C(a,b)} \int_0^1 \mu^a (1-\mu)^{b-1} d\mu. \end{aligned} \quad (24)$$

The integration term of Equation 24,  $\int_0^1 \mu^a (1-\mu)^{b-1} d\mu$ , also follows the cumulative distribution formulation of the Beta distribution. Therefore, we can further re-write Equation 24 term as,

$$\begin{aligned} \int_0^1 \mu^a (1-\mu)^{b-1} d\mu &= \int_0^1 \mu^{(a+1)-1} (1-\mu)^{b-1} d\mu \\ &= \int_0^1 \mu^{(a+1)-1} (1-\mu)^{b-1} d\mu \\ &= C(a+1, b) \\ &= \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)}. \end{aligned} \quad (25)$$

By plugging the above term in Equation 24, we can re-write  $E(\mu)$ , as

$$\begin{aligned} E(\mu) &= \frac{1}{C(a,b)} \int_0^1 \mu^a (1-\mu)^{b-1} d\mu \\ &= \frac{1}{C(a,b)} \times \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \times \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \times \frac{a\Gamma(a)\Gamma(b)}{(a+b)\Gamma(a+b)} \\ &= \frac{a}{a+b}. \end{aligned} \quad (26)$$

#### IV. SUPPORT VECTOR MACHINES

Recall that the goal of a support vector machine (SVM) classifier is maximizing the geometric margin  $\rho$ , given data  $\mathbf{x}_i$  with labels  $y_i$ ,

$$\max_{\mathbf{w}, b} \rho = \min_i \frac{y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)}{\|\mathbf{w}\|}. \quad (27)$$

Here,  $i \in \{1, 2, \dots, N\}$ ,  $N$  is the number of data, and  $\phi(\cdot)$  is a kernel function that maps  $\mathbf{x}_i$  into a hyperspace of more dimensions, and  $b$  is the bias. By scaling as  $\mathbf{w} \rightarrow \kappa \mathbf{w}$  and  $b \rightarrow \kappa b$ , we can get

$$\min_i y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) = 1, \quad (28)$$

and re-formulate the problem as,

$$\max_{\mathbf{w}, b} \rho = \frac{1}{\|\mathbf{w}\|} \quad (29)$$

$$\text{s.t. } y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) \leq 1. \quad (30)$$

Then, this optimization can be transformed as

$$\max_{\mathbf{w}, b} \rho \Leftrightarrow \min_{\mathbf{w}, b} \frac{1}{\rho^2} \Leftrightarrow \max_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (31)$$

where the same constraint of Equation 30 is given.

By using the method of Lagrange multipliers, Equation 30, 31 are formulated, as

$$\min_{\mathbf{w}, b, \boldsymbol{\lambda}} L(\mathbf{w}, b, \boldsymbol{\lambda}) \quad (32)$$

where

$$L(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \lambda_i (y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) - 1). \quad (33)$$

The solutions of the problem are computed from the partial derivatives,

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^N \lambda_i y_i \phi(\mathbf{x}_i) = 0, \quad (34)$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^N \lambda_i y_i = 0. \quad (35)$$

Also, a constrained optimization must satisfy the Karush-Kuhn-Tucker (KKT) conditions [1], which yields the conditions,

$$\begin{aligned} \forall i, \quad \lambda_i &> 0, \\ \lambda_i (y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) - 1) &= 0. \end{aligned} \quad (36)$$

By using Equation 36, we can get the following equation,

$$\sum_{i=1}^N \lambda_i (y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) - 1) = 0. \quad (37)$$

This yields

$$\begin{aligned} 0 &= \sum_{i=1}^N \lambda_i y_i \mathbf{w}^\top \phi(\mathbf{x}_i) + \sum_{i=1}^N \lambda_i y_i b - \sum_{i=1}^N \lambda_i \\ &= \mathbf{w}^\top \sum_{i=1}^N \lambda_i y_i \phi(\mathbf{x}_i) + b \sum_{i=1}^N \lambda_i y_i - \sum_{i=1}^N \lambda_i. \end{aligned} \quad (38)$$

By plugging Equation 34, 35 in the above equation, we can get

$$\begin{aligned}
 0 &= \mathbf{w}^\top \sum_{i=1}^N \lambda_i y_i \phi(\mathbf{x}_i) + b \sum_{i=1}^N \lambda_i y_i - \sum_{i=1}^N \lambda_i = \mathbf{w}^\top \mathbf{w} + b \times 0 - \sum_{i=1}^N \lambda_i \\
 &= \|\mathbf{w}\|^2 - \sum_{i=1}^N \lambda_i.
 \end{aligned} \tag{39}$$

Thus, we can show that

$$\|\mathbf{w}\|^2 = \sum_{i=1}^N \lambda_i. \tag{40}$$

## V. GIBBS SAMPLING

The Gaussian random variable  $\nu$  is independent from  $z_i$ , and  $E(\nu) = 0$ . Therefore, we can get

$$E(\sigma_i \nu) = E(\sigma_i)E(\nu) = E(\sigma_i) \times 0 = 0. \quad (41)$$

Therefore, the expectation of  $\hat{z}_i$  can be re-written, as

$$\begin{aligned} E(\hat{z}_i) &= E\left(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha^2)^{\frac{1}{2}}\nu\right) \\ &= (1 - \alpha)E(\mu_i) + \alpha E(z_i) + (1 - \alpha^2)^{\frac{1}{2}}E(\sigma_i \nu) \\ &= (1 - \alpha)E(\mu_i) + \alpha E(z_i) + (1 - \alpha^2)^{\frac{1}{2}} \times 0 \\ &= (1 - \alpha)E(\mu_i) + \alpha E(z_i). \end{aligned} \quad (42)$$

From the definition,  $\mu_i = E(z_i)$ . Also,  $E(\mu_i) = \mu_i$ . Thus, we can further simplify Equation 42, as

$$\begin{aligned} E(\hat{z}_i) &= (1 - \alpha)\mu_i + \alpha E(z_i) \\ &= (1 - \alpha)\mu_i + \alpha \mu_i \\ &= \mu_i. \end{aligned} \quad (43)$$

Therefore, the mean of  $\hat{z}_i$  is also  $\mu_i$ .

## REFERENCES

- [1] Christopher M. Bishop (2006) *Pattern Recognition and Machine Learning*, Springer.