

# CS405 Homework #1

12110644 周思呈

## Question 1

Consider the polynomial function:

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{i=0}^M w_i x^i$$

Calculate the coefficients  $\mathbf{w} = w_i$  that minimize its sum-of-squares error function. Here a suffix  $i$  denotes the index of a component, whereas  $(x)^i$  denotes  $x$  raised to the power of  $i$ .

Let  $\mathbf{x} = [x^0, x^1, \dots, x^m]^T$ ,  $\mathbf{w} = [w_0, w_1, \dots, w_M]$  so  $y(x, \mathbf{w}) = \mathbf{w}\mathbf{x} = \mathbf{x}^T \mathbf{w}^T$ .

Assume our target label value is  $t_n$ , so we have to minimize

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \quad (1)$$

where we add coefficient  $\frac{1}{2}$  to simplify the derivation process. Let  $\mathbf{t} = [t_0, t_1, \dots, t_n]$ ,  $\mathbf{X} = [\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n]$ .

$$\begin{aligned} \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} &= \sum_{n=1}^N (y_n \frac{\partial y}{\partial \mathbf{w}} - t_n \frac{\partial y}{\partial \mathbf{w}}) \\ &= \sum_{n=1}^N \frac{\partial (\mathbf{x}^T \mathbf{w}^T \mathbf{w} \mathbf{x} - 2t_n \mathbf{w} \mathbf{x})}{\partial \mathbf{w}} \\ &= \sum_{n=1}^N (2\mathbf{w} \mathbf{x} \mathbf{x}^T - 2t_n \mathbf{x}^T) \\ &= 2(\mathbf{w} \mathbf{X} \mathbf{X}^T - \mathbf{t} \mathbf{X}^T) = 0 \end{aligned} \quad (2)$$

So  $\mathbf{w} = \mathbf{t} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1}$ .

## Question 2

Suppose that we have three colored boxes  $r(\text{red})$ ,  $b(\text{blue})$ , and  $g(\text{green})$ . Box  $r$  contains 3 apples, 4 oranges, and 3 limes, box  $b$  contains 1 apple, 1 orange, and 0 limes, and box  $g$  contains 3 apples, 3 oranges, and 4 limes. If a box is chosen at random with probabilities  $p(r) = 0.2, p(b) = 0.2, p(g) = 0.6$ , and a piece of fruit is removed from the box (with equal probability of selecting any of the items in the box), then what is the probability of selecting an apple? If we observe that the selected fruit is in fact an orange, what is the probability that it came from the green box?

Let  $p(a)$  indicates selecting an apple, then

$$\begin{aligned} p(a) &= p(a|r)p(r) + p(a|g)p(g) + p(a|b)p(b) \\ &= \frac{3}{3+4+3} \times 0.2 + \frac{3}{3+3+4} \times 0.6 + \frac{1}{1+1} \times 0.2 \\ &= 0.34 \end{aligned} \quad (3)$$

So the probability of selecting an apple is 0.34.

Let  $p(o)$  indicates selecting an orange, then according to the Bayes' Law,

$$\begin{aligned} p(g|o) &= \frac{p(g)p(o|g)}{p(o)} \\ &= \frac{0.6 \times 0.3}{0.36} \\ &= \frac{1}{2} \end{aligned} \quad (4)$$

If the selected fruit is an orange, the possibility it came from the green box is 0.5.

## Question 3

Given two statistically independent variables  $x$  and  $z$ , show that the mean and variance of their sum satisfies

$$\begin{aligned}\mathbb{E}[x + z] &= \mathbb{E}[x] + \mathbb{E}[z] \\ \text{var}[x + z] &= \text{var}[x] + \text{var}[z]\end{aligned}\tag{5}$$

Suppose  $X$  and  $Z$  are continuous.

$$\begin{aligned}E(X + Z) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x + z)p(x, z)dx dz \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xp(x, z)dx dz + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} zp(x, z)dx dz \\ &= \int_{-\infty}^{\infty} xp_X(x)dx + \int_{-\infty}^{\infty} zp_Z(z)dz \\ &= E(X) + E(Z)\end{aligned}\tag{6}$$

It is similar in the discrete case.

$$\begin{aligned}E[X + Z] &= \sum_{i,j} (x_i + z_j)P\{X = x_i, Z = z_j\} \\ &= \sum_{i,j} x_i P\{X = x_i, Z = z_j\} + \sum_{i,j} z_j P\{X = x_i, Z = z_j\} \\ &= \sum_i x_i \sum_j P\{X = x_i, Z = z_j\} + \sum_j z_j \sum_i P\{X = x_i, Z = z_j\} \\ &= \sum_i x_i P\{X = x_i\} + \sum_j z_j P\{Z = z_j\} \\ &= E[X] + E[Z]\end{aligned}\tag{7}$$

Variance of variable  $X$  is defined as

$$D(X) = E((X - E(X))^2)\tag{8}$$

And it has a property that

$$\begin{aligned}D(X) &= E\{[X - E(X)]^2\} = E\{X^2 - 2XE(X) + [E(X)]^2\} \\ &= E(X^2) - 2E(X)E(X) + [E(X)]^2 \\ &= E(X^2) - [E(X)]^2\end{aligned}\tag{9}$$

Apply it to what we need to prove

$$\begin{aligned}D(X + Z) &= E[(X + Z)^2] - E^2(X + Z) \\ &= E(X^2 + 2XZ + Z^2) - (E(X) + E(Z))^2 \\ &= E(X^2) + 2E(XZ) + E(Z^2) - E^2(X) - 2E(X)E(Z) - E^2(Z) \\ &= D(X) + D(Z) + 2Cov(X, Z)\end{aligned}\tag{10}$$

Because  $X$  and  $Z$  are independent, we have  $Cov(X, Z) = 0$ , so  $D(X + Z) = D(X) + D(Z)$ .

## Question 4

In probability theory and statistics, the Poisson distribution, is a discrete probability distribution that expresses the probability of a given number of events occurring in a fixed interval of time or space if these events occur with a known constant rate and independently of the time since the last event. If  $X$  is Poisson distributed, i.e.  $X \sim \text{Poisson}(\lambda)$ , its probability mass function takes the following form:

$$P(X|\lambda) = \frac{\lambda^X e^{-\lambda}}{X!}$$

It can be shown that if  $\mathbb{E}(X) = \lambda$ . Assume now we have  $n$  data points from  $\text{Poisson}(\lambda) : \mathcal{D} = \{X_1, X_2, \dots, X_n\}$ . Show that the sample mean  $\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n X_i$  is the maximum likelihood estimate(MLE) of  $\lambda$ .

If  $X$  is exponential distribution and its distribution density function is  $f(x) = \frac{1}{\lambda} e^{-\frac{x}{\lambda}}$  for  $x > 0$  and  $f(x) = 0$  for  $x \leq 0$ . Show that the sample mean  $\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n X_i$  is the maximum likelihood estimate(MLE) of  $\lambda$ .

The likelihood function is

$$lik(\theta) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) \quad (11)$$

If the random variables are independently and identically distributed(IID), we have

$$lik(\theta) = P(X_1 = x_1)P(X_2 = x_2) \dots P(X_n = x_n) \quad (12)$$

Consider random variables in Poisson distribution

$$L(x|\lambda) = \frac{\lambda^{\sum_{i=1}^n x_i} e^{-n\lambda}}{\prod_{i=1}^n x_i!} \quad (13)$$

We want to find out what will  $\lambda$  be when  $L$  is the smallest, so we take the log of  $L$  and take the partial derivative with respect to  $\lambda$

$$\frac{\partial L(x|\lambda)}{\partial \lambda} = \sum_{i=1}^n x_i \frac{1}{\lambda} - n = 0 \quad (14)$$

So we have  $\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n X_i$ .

If  $X$  is exponential distribution, consider when  $x > 0$

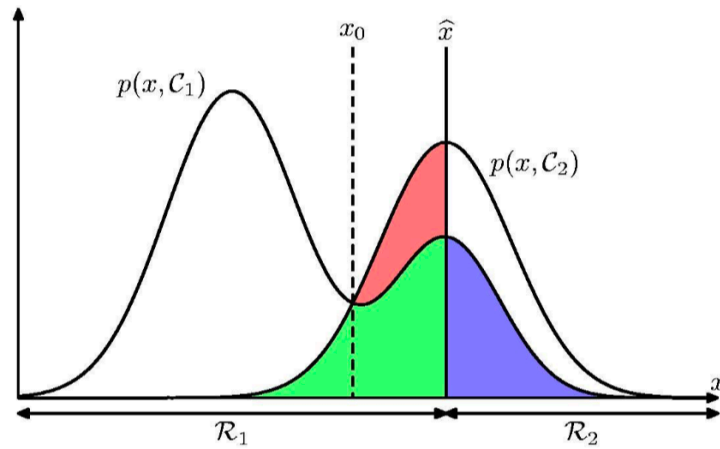
$$L(x|\lambda) = \left(\frac{1}{\lambda}\right)^n e^{-\frac{1}{\lambda} \sum_{i=1}^n x_i} \quad (15)$$

$$\frac{\partial L(x|\lambda)}{\partial \lambda} = -n\lambda^{-1} + \sum_{i=1}^n \lambda^{-2} = 0$$

So we have  $\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n X_i$ .

## Question 5

(a) Write down the probability of classifying correctly  $p(\text{correct})$  and the probability of misclassification  $p(\text{mistake})$  according to the following chart.



$$p(\text{mistake}) = p(x \in \mathcal{R}_1, C_2) + p(x \in \mathcal{R}_2, C_1)$$

$$= \int_{\mathcal{R}_1} p(x, C_2) dx + \int_{\mathcal{R}_2} p(x, C_1) dx \quad (16)$$

$$p(\text{correct}) = \sum_{k=1}^K p(x \in \mathcal{R}_k, C_k) = \sum_{k=1}^K \int_{\mathcal{R}_k} p(x, C_k) dx \quad (17)$$

(b) For multiple target variables described by vector  $\mathbf{t}$ , the expected squared loss function is given by

$$\mathbb{E}[L(\mathbf{t}, \mathbf{y}(\mathbf{x}))] = \int \int \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$$

Show that the function  $\mathbf{y}(\mathbf{x})$  for which this expected loss is minimized given by  $\mathbf{y}(\mathbf{x}) = \mathbb{E}_{\mathbf{t}}[\mathbf{t}|\mathbf{x}]$ .

## Hints

For a single target variable  $t$ , the loss is given by

$$\mathbb{E}[L] = \int \int \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

The result is as follows

$$y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int t p(t|\mathbf{x}) dt = \mathbb{E}_t[t|\mathbf{x}]$$

$$\begin{aligned} \frac{\delta \mathbb{E}(L)}{\delta \mathbf{y}(\mathbf{x})} &= 2 \left[ \int (y_1(\mathbf{x}) - t) p(\mathbf{x}, t) dt, \dots, \int (y_n(\mathbf{x}) - t) p(\mathbf{x}, t) dt \right] = 0 \\ y_i(\mathbf{x}) &= \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = E[t|\mathbf{x}] \end{aligned} \quad (18)$$

So we have  $\mathbf{y}(\mathbf{x}) = \mathbb{E}_t[\mathbf{t}|\mathbf{x}]$ .

## Question 6

(a) We defined the entropy based on a discrete random variable  $\mathbf{X}$  as

$$\mathbf{H}[\mathbf{X}] = - \sum_i p(x_i) \ln p(x_i)$$

Now consider the case that  $\mathbf{X}$  is a continuous random variable with the probability density function  $p(x)$ . The entropy is defined as

$$\mathbf{H}[\mathbf{X}] = - \int p(x) \ln p(x) dx$$

Assume that  $\mathbf{X}$  follows Gaussian distribution with the mean  $\mu$  and variance  $\sigma$ , i.e.

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Please derive its entropy  $\mathbf{H}[\mathbf{X}]$ .

$$\begin{aligned} H[x] &= - \int \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\} \ln \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\} dx \\ &= - \frac{1}{(2\pi\sigma^2)^{1/2}} \int \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\} \left( -\ln(\sqrt{2\pi}\sigma) - \frac{(x-\mu)^2}{2\sigma^2} \right) dx \\ &= - \frac{1}{(2\pi\sigma^2)^{1/2}} \cdot -\ln(\sqrt{2\pi}\sigma) \int \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\} dx + \frac{1}{(2\pi\sigma^2)^{1/2}} \int \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\} \frac{(x-\mu)^2}{2\sigma^2} dx \\ &= \frac{\ln(\sqrt{2\pi}\sigma)}{(2\pi\sigma^2)^{1/2}} \int \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\} dx + \frac{1}{(2\pi\sigma^2)^{1/2}} \int \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\} \frac{(x-\mu)^2}{2\sigma^2} dx \\ &= \frac{\ln(\sqrt{2\pi}\sigma)}{(2\pi\sigma^2)^{1/2}} \sqrt{2\sigma} \int \exp \left\{ -\left( \frac{x-\mu}{\sqrt{2}\sigma} \right)^2 \right\} d \left( \frac{x-\mu}{\sqrt{2}\sigma} \right) + \frac{1}{(2\pi\sigma^2)^{1/2}} \sqrt{2\sigma} \int \exp \left\{ -\left( \frac{x-\mu}{\sqrt{2}\sigma} \right)^2 \right\} \frac{(x-\mu)^2}{2\sigma^2} d \left( \frac{x-\mu}{\sqrt{2}\sigma} \right) \\ &= \frac{\ln(\sqrt{2\pi}\sigma)}{\sqrt{\pi}} \int_{-\infty}^{\infty} e^{-y^2} dy + \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} e^{-y^2} y^2 dy \\ &= \ln(\sqrt{2\pi}\sigma) + \frac{1}{\sqrt{\pi}} \cdot -\frac{1}{2} \left( 0 - \int_{-\infty}^{\infty} e^{-y^2} dy \right) \\ &= \ln(\sqrt{2\pi}\sigma) + \frac{1}{2} \\ &= \frac{1}{2} (\ln(2\pi\sigma^2) + 1) \end{aligned}$$

(b) Write down the mutual information  $\mathbf{I}(\mathbf{y}, \mathbf{x})$ . Then show the following equation

$$\mathbf{I}[\mathbf{x}, \mathbf{y}] = \mathbf{H}[\mathbf{x}] - \mathbf{H}[\mathbf{x}|\mathbf{y}] = \mathbf{H}[\mathbf{y}] - \mathbf{H}[\mathbf{y}|\mathbf{x}]$$

Mutual information is the relative entropy of the joint distribution and the product distribution. Consider the continuous situation

$$\begin{aligned}
\mathbf{I}[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) \\
&= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left( \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x}d\mathbf{y} \\
\mathbf{H}[\mathbf{x}] &= - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \\
\mathbf{H}[\mathbf{y}] &= - \int p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{y} \\
\mathbf{I}[\mathbf{x}, \mathbf{y}] &= \iint p(\mathbf{x}, \mathbf{y}) \ln \left( \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})} \right) d\mathbf{x}d\mathbf{y} - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \\
&= \iint p(\mathbf{y})p(\mathbf{x}|\mathbf{y}) \ln (p(\mathbf{x}|\mathbf{y})) d\mathbf{x}d\mathbf{y} + \mathbf{H}[\mathbf{x}] \\
&= -\mathbf{H}[\mathbf{x}|\mathbf{y}] + \mathbf{H}[\mathbf{x}]
\end{aligned} \tag{19}$$

By analogy, we have  $\mathbf{I}[\mathbf{x}, \mathbf{y}] = \mathbf{H}[\mathbf{y}] - \mathbf{H}[\mathbf{y}|\mathbf{x}]$ .