# CS329 Machine Learning

## Homework #1

**Fan Site**
fanst2021@mail.sustech.edu.cn

# Question 1

Consider the polynomial function:

$$y(x, \boldsymbol{w}) = w_0 + w_1 x + w_2 x + \ldots + w_M x^M = \sum_{i=0}^{M} w_i x^i$$

Calculate the coefficients $\boldsymbol{w} = \{w_i\}$ 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 $i$.

> ## Solution: least squares
>
> The sum-of-squares error function is defined as:
>
> $$E(\boldsymbol{w}) = \sum_{j=1}^{N} \left[ y(x_j, \boldsymbol{w}) - t_j \right]^2$$
>
> where:
> - $N$ is the number of data points.
> - $(x_j, t_j)$ are the input-output pairs in the dataset.
> - $y(x_j, \boldsymbol{w})$ is the given polynomial function.
>
> To minimize the error function $E(\boldsymbol{w})$ with respect to the coefficients $\boldsymbol{w} = \{w_i\}$, we'll take the derivative of the error function with respect to each coefficient $w_i$ and set each derivative equal to zero to find the minimum. The derivative of $E(\boldsymbol{w})$ with respect to $w_i$ is:
>
> $$\frac{\partial E}{\partial w_i} = -2 \sum_{j=1}^{N} x_j^i \left( t_j - \sum_{k=0}^{M} w_k x_j^k \right)$$
>
> Setting this derivative equal to zero for each $i$:
>
> $$-2 \sum_{j=1}^{N} x_j^i \left( t_j - \sum_{k=0}^{M} w_k x_j^k \right) = 0$$
>
> Now, we can solve for the coefficients $w_i$:
>
> $$\sum_{j=1}^{N} x_j^i \sum_{k=0}^{M} w_k x_j^k = \sum_{j=1}^{N} x_j^i t_j$$
>
> This equation forms a system of linear equations for the coefficients $\boldsymbol{w}$.
>
> $$\boldsymbol{X} \boldsymbol{w} = \boldsymbol{T}$$
>
> Where:
> - $\boldsymbol{X}$ is the design matrix whose elements are $x_j^i$ for all $i$ and $j$.
> - $\boldsymbol{T}$ is the target vector whose elements are the corresponding $t_j$ values.
>
> Hence $\boldsymbol{w} = \left( \boldsymbol{X}^T \boldsymbol{X} \right)^{-1} \boldsymbol{X}^T \boldsymbol{T}$
>
> With this $\boldsymbol{w}$, the sum-of-squares error is minimized.

# Question 2

Suppose that we have three colored boxes $r$(red), $b$(blue), and $g$(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?

## Question 3

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

$$\mathbb{E}[x + z] = \mathbb{E}[x] + \mathbb{E}[z]$$

$$\text{var}[x + z] = \text{var}[x] + \text{var}[z]$$

$$\mathbb{E}(X+Z) = \sum_x \sum_z (x+z)p(x,z)$$

$$= \sum_x \sum_z xp(x,z) + \sum_x \sum_z zp(x,z)$$

$$= \sum_x xp_X(x) + \sum_z xp_Z(z)$$

$$= \mathbb{E}(X) + \mathbb{E}(Z)$$

Then prove that for independent variables $X$ and $Z$, "the variance of the sum equals to sum of the variance"

$$\text{var}[X+Z] = \mathbb{E}\big((X+Z)^2\big) - (\mathbb{E}(X+Z))^2$$

$$= \mathbb{E}(X^2 + 2XZ + Z^2) - \big((\mathbb{E}(X))^2 + 2\mathbb{E}(X)\mathbb{E}(Z) + (\mathbb{E}(Z))^2\big)$$

$$= \mathbb{E}(X^2) - (\mathbb{E}(X))^2 + \mathbb{E}(Z^2) - (\mathbb{E}(Z))^2 + 2(\mathbb{E}(XZ) - \mathbb{E}(X)\mathbb{E}(Z))$$

$$= \text{var}[X] + \text{var}[Z]$$

# 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. $X \sim \text{Possion}(\lambda)$, 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 $\text{Possion}(\lambda) : \mathcal{D} = \{X_1, X_2, ..., X_n\}$. 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 $\lambda$.

## Solution

**Possion distribution**

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

$$\text{Likelihood} \quad L(\lambda) = \prod_{i=1}^{n} P(X_i|\lambda)$$

$$\ln L(\lambda) = \ln\left(\prod_{i=1}^{n} \frac{\lambda_i^X e^{-\lambda}}{X_i!}\right) = \ln\left(\prod_{i=1}^{n} \frac{1}{X_i!}\right) + \ln\lambda \sum_{i=1}^{n} X_i - n\lambda$$

$$\text{Let} \quad \frac{\partial \ln L(\lambda)}{\partial \lambda} = \sum_{i=1}^{n} \frac{X_i}{\lambda} - n = 0$$

$$\text{The MLE is } \hat{\lambda} = \frac{1}{n}\sum_{i=1}^{n} X_i.$$

**Exponential distribution**

$$P(X|\lambda) = \frac{1}{\lambda} e^{-\frac{X}{\lambda}}$$

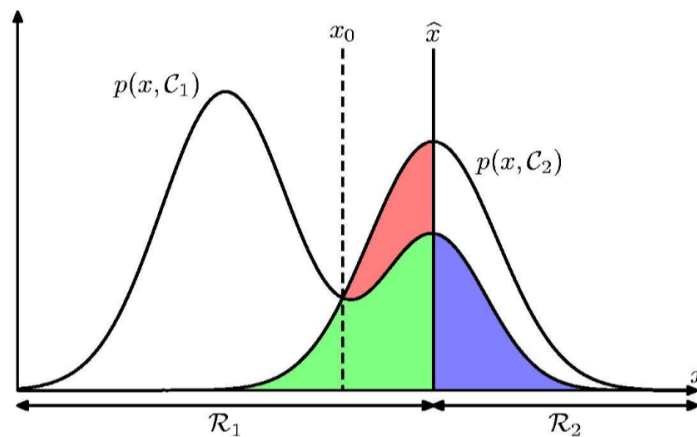$$\text{Likelihood} \quad L(\lambda) = \prod_{i=1}^{n} P(X_i|\lambda)$$

$$\ln L(\lambda) = \ln\left(\prod_{i=1}^{n} \frac{1}{\lambda} e^{-\frac{X_i}{\lambda}}\right) = -n\ln\lambda - \sum_{i=1}^{n} \frac{X_i}{\lambda}$$

$$\text{Let} \quad \frac{\partial \ln L(\lambda)}{\partial \lambda} = -\frac{n}{\lambda} + \sum_{i=1}^{n} \frac{X_i}{\lambda^2} = 0$$

$$\text{The MLE is } \hat{\lambda} = \frac{1}{n}\sum_{i=1}^{n} X_i.$$

# 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.



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

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

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

### (a) Solution

$$p(\text{correct}) = p(x \in \mathcal{R}_1, \mathcal{C}_1) + p(x \in \mathcal{R}_2, \mathcal{C}_2) = \int_{\mathcal{R}_1} p(x, \mathcal{C}_1)dx + \int_{\mathcal{R}_2} p(x, \mathcal{C}_2)dx$$

$$p(\text{mistake}) = p(x \in \mathcal{R}_1, \mathcal{C}_2) + p(x \in \mathcal{R}_2, \mathcal{C}_1) = \int_{\mathcal{R}_1} p(x, \mathcal{C}_2)dx + \int_{\mathcal{R}_2} p(x, \mathcal{C}_1)dx$$

### (b) Solution

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

$$\text{Let} \quad \frac{\partial \mathbb{E}[L(t, y(x))]}{\partial y(x)} = 2\int (y(x) - t)p(x, t)dt = 0$$

$$\text{Isolating } y(x), \quad y(x)\int p(x, t)dt = \int tp(x, t)dt$$

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

## Question 6

(a) We defined the entropy based on a discrete random variable $X$ as

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

Now consider the case that $X$ 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 $X$ 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 $\mathbf{H}[\mathbf{X}]$.

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

$$\mathbf{I}[\mathbf{X},\mathbf{Y}] = \mathbf{H}[\mathbf{X}] - \mathbf{H}[\mathbf{X}|\mathbf{Y}] = \mathbf{H}[\mathbf{Y}] - \mathbf{H}[\mathbf{Y}|\mathbf{X}]$$

### (a) Solution

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

$$\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$$

$$\text{Let} \quad t = \frac{x-\mu}{\sqrt{2}\sigma}, \frac{dt}{dx} = \frac{1}{\sqrt{2}\sigma}$$

$$\mathbf{H[X]} = -\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \ln\left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}\right) dx$$

$$= -\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-t^2} \ln\left(\frac{1}{\sqrt{2\pi}\sigma} e^{-t^2}\right) dx$$

$$= -\frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{-t^2} \left(-\ln\left(\sqrt{2\pi}\sigma\right) - t^2\right) dx$$

$$= \frac{\ln\left(\sqrt{2\pi}\sigma\right)}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{-t^2} dx + \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} t^2 e^{-t^2} dx$$

$$= \frac{\ln\left(\sqrt{2\pi}\sigma\right)}{\sqrt{\pi}} \int_{-\infty}^{\infty} e^{-t^2} dt + \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} t^2 e^{-t^2} dt$$

$$= \ln\left(\sqrt{2\pi}\sigma\right) + \frac{1}{\sqrt{\pi}} \cdot -\frac{1}{2}\left(0 - \int_{-\infty}^{\infty} e^{-t^2} dt\right)$$

$$= \ln\left(\sqrt{2\pi}\sigma\right) + \frac{1}{2}$$

## (b) Solution

**Discrete Distributions**

$$\mathbf{I[X,Y]} = \sum_{x\in\mathcal{X}}\sum_{y\in\mathcal{Y}} P_{X,Y}(x,y) \ln\left(\frac{P_{X,Y}(x,y)}{P_X(x)P_Y(y)}\right)$$

$$\mathbf{I[X,Y]} = \sum_{x\in\mathcal{X}}\sum_{y\in\mathcal{Y}} P_{X,Y}(x,y) \ln\left(\frac{P_{X,Y}(x,y)}{P_Y(y)}\right) - \sum_{x\in\mathcal{X}} P_X(x) \ln(P_X(x))$$

$$= \sum_{x\in\mathcal{X}}\sum_{y\in\mathcal{Y}} P_Y(y) P(x|y) \ln P(x|y) + \mathbf{H(X)}$$

$$= \mathbf{H[X]} - \mathbf{H[X|Y]}.$$

where $P_{X,Y}$ is the joint probability mass function of $X$ and $Y$, and $P_X$ and $P_Y$ are the marginal probability mass functions of $X$ and $Y$ respectively. Similarly, $\mathbf{I[X,Y]} = \mathbf{H[Y]} - \mathbf{H[Y|X]}$.

**Continuous Distributions**

$$\mathbf{I[x,y]} = \int_{\mathcal{X}}\int_{\mathcal{Y}} P_{X,Y}(x,y) \ln\left(\frac{P_{X,Y}(x,y)}{P_X(x)P_Y(y)}\right) dy dx$$

$$\mathbf{I[X,Y]} = \int_{\mathcal{X}}\int_{\mathcal{Y}} P_{X,Y}(x,y) \ln\left(\frac{P_{X,Y}(x,y)}{P_Y(y)}\right) dy dx - \int_{\mathcal{X}} P_X(x) \ln(P_X(x)) dx$$

$$= \int_{\mathcal{X}}\int_{\mathcal{Y}} P_Y(y) P(x|y) \ln P(x|y) dy dx + \mathbf{H(X)}$$

$$= \mathbf{H[X]} - \mathbf{H[X|Y]}.$$

where $P_{X,Y}$ is the joint probability density function of $X$ and $Y$, and $P_X$ and $P_Y$ are the marginal probability density functions of $X$ and $Y$ respectively. Similarly, $\mathbf{I[X,Y]} = \mathbf{H[Y]} - \mathbf{H[Y|X]}$.