# HOMEWORK 4

| Full name | Nguyễn Thị Minh Ngọc |
|---|---|
| Student ID | 11219280 |
| Class | DSEB 63 |

## 1.

We have:
- The a data set of observations: $\mathbf{x} = (x_1, x_2 \dots x_N)^T$ , representing N observations of the scalar variable x
- The corresponding target values of $\mathbf{x}$: $\mathbf{t} = (x_1, x_2 \dots x_N)^T$

Based on those data, we need to find a model which can help us to make predictions for some new value of the input variable x.

Suppose that the observations are drawn independently from a Gaussian distribution.
$$t = y(x, \mathbf{w}) + noise$$
With noise represents factors that affect the output but cannot be evaluated easily. Suppose that noise $\sim N(0, \beta^{-1}) \rightarrow t = y(x, w) + noise \sim N(y(x, \mathbf{w}), \beta^{-1})$ where $\beta = \frac{1}{\sigma^2}$. Then:
$$p(t|x, \mathbf{w}, \beta^{-1}) = N(t|x, \mathbf{w}, \beta^{-1})$$
To find the best model, we need to determine the unknown parameter $\mathbf{w}$ to maximize $p(t|x, \mathbf{w}, \beta^{-1})$. The likelihood function:
$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta^{-1}) = \prod_{n=1}^{N} N(t|x, \mathbf{w}, \beta^{-1})$$
Because $0 < p(t|x, \mathbf{w}, \beta^{-1}) < 1$, so the product of N p(t) will come to almost 0. Hence, it is convenient to maximize the logarithm of the likelihood function:

| $\log\big(p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta^{-1})\big)$ | $= \sum_{n=1}^{N} \log\left( \frac{1}{\sqrt{2\pi^{\beta^{-1}}}} e^{\frac{-(t-y(x,\mathbf{w}))^2 \beta}{2}} \right)$ |
|---|---|
| | $= N\log\frac{1}{\sqrt{2\pi^{\beta^{-1}}}} + \sum_{n=1}^{N} \log\left( e^{\frac{-(t-y(x,\mathbf{w}))^2\beta}{2}} \right)$ |
| | $= N\log\frac{1}{\sqrt{2\pi^{\beta^{-1}}}} + \sum_{n=1}^{N} \log\left( e^{\frac{-(t-y(x,\mathbf{w}))^2\beta}{2}} \right)$ |
| | $= N\log\frac{1}{\sqrt{2\pi^{\beta^{-1}}}} - \frac{\beta}{2}\sum_{n=1}^{N} (y(x, \mathbf{w}) - t)^2$ |

To maximize $\log\left(p(\mathbf{t}|\mathbf{x},\mathbf{w},\beta^{-1})\right)$ and find $\mathbf{w}$, we need to minimize $\frac{\beta}{2}\sum_{n=1}(y(x,w)-t)^2$ or minimize $\sum_{n=1}(y(x,\mathbf{w})-t)^2$.

Let $P = \sum_{n=1}(y(x,\mathbf{w})-t)^2$ where $y = w_1 x + w_0$

Suppose that:

$$X = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix}; \qquad \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}; \qquad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} w_1 x_1 + w_0 \\ w_2 x_2 + w_0 \\ \vdots \\ w_n x_n + w_0 \end{bmatrix} = X \cdot \mathbf{w}$$

Then $P = \|y(x_1\mathbf{w})-t\|_2^2 = \|X.\mathbf{w}-t\|_2^2$

We have:

$$\frac{\partial P}{\partial \mathbf{w}} = \begin{bmatrix} \dfrac{\partial P}{\partial w_0} \\ \dfrac{\partial P}{\partial w_1} \end{bmatrix} = \begin{bmatrix} 2(X\mathbf{w}-t) \\ 2X(X\mathbf{w}-t) \end{bmatrix} = 2X^T(X\mathbf{w}-t) = 0$$

$$\to 2X^T\mathbf{w}X - 2X^T t = 0$$
$$\to X^T\mathbf{w}X = X^T t$$
$$\to \mathbf{w} = (X^{\wedge}TX)^{-1}X^T t$$

**Extra:** Prove that X^TX is invertible when X is full of rank

…