

- ③ The method of ordinary least squares assumes that there is constant variance in the errors.

Let the data observed $D = \{(x_1, t_1), (x_2, t_2), \dots, (x_n, t_n)\}$

Let us assume that the targets are generated by $y(x, w)$

$$\text{i.e., } t = y(x, w) + \varepsilon \quad \varepsilon \sim N(\varepsilon | 0, \sigma^2)$$

where ε follows a normal distribution with varying variance.

- a) likelihood for a heteroscedastic setting of a single data point =

$$p(t_n | x_n, w, \sigma_n^2) = \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{1}{2\sigma_n^2} (t_n - y(x_n, w))^2\right)$$

Now, we are going to model the prior distribution also with normal distribution.

prior =

$$p(w | \alpha) = N(w_n | 0, \alpha^{-1}) \quad \alpha = \text{precision parameter}$$

$$= \left(\frac{\alpha}{2\pi}\right)^{1/2} \cdot \exp\left(-\frac{\alpha}{2} w_n \cdot w_n\right)$$

b) Maximum likelihood estimation:

$$\begin{aligned} \text{Likelihood} &= p(t_1, \dots, t_n | x_1, \dots, x_n, \omega, \sigma^2) \\ &= p(t_1 | x_1, \omega, \sigma_1^2) \cdot p(t_2 | x_2, \omega, \sigma_2^2) \cdot \dots \cdot p(t_n | x_n, \omega, \sigma_n^2) \end{aligned}$$

$$\begin{aligned} &= \prod_{i=1}^n p(t_i | x_i, \omega, \sigma_i^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(t_i - y(x_i, \omega))^2}{2\sigma_i^2}} \end{aligned}$$

Instead of working with likelihood, it will be easier to work with loglikelihood.

$$\log(\text{likelihood}) = \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi\sigma_i^2}}\right) - \sum_{i=1}^n \frac{1}{2\sigma_i^2} (t_i - y(x_i, \omega))^2$$

This is the objective function that will be considered for the ML estimate of parameters.

$$\Rightarrow \omega_{ML} = \arg \max_{\omega} \left[\underbrace{\sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi\sigma_i^2}}\right)}_{\substack{\downarrow \\ \text{this does not} \\ \text{depend on } \omega}} - \sum_{i=1}^n \frac{1}{2\sigma_i^2} (t_i - y(x_i, \omega))^2 \right]$$

$$\Rightarrow \omega_{ML} = \arg \min_{\omega} \frac{1}{2} \sum_{i=1}^n \frac{1}{\sigma_i^2} (t_i - y(x_i, \omega))^2$$

Ma

MAP estimation :

Given a prior, the posterior distribution is given by:

$$p(w|x, t, \beta, \alpha) = \frac{p(t|x, w, \beta) \cdot p(w|\alpha)}{p(t|x, \beta, \alpha)}$$

we have to maximize this

this does not depend on w

$$p(w|\alpha) = \prod_{i=1}^N N(w_i | 0, \alpha^{-1}) = \prod_{i=1}^N \left(\frac{\alpha}{2\pi} \right)^{1/2} e^{-\frac{\alpha}{2} w_i^2}$$

$$= \left(\frac{\alpha}{2\pi} \right)^{N/2} \prod_{i=1}^N e^{-\frac{\alpha}{2} w_i^2}$$

$$\log p(w|\alpha) = \log \left(\frac{\alpha}{2\pi} \right)^{N/2} - \frac{\alpha}{2} w^T w$$

$$p(t|x, w, \beta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(t_i - y(x_i, w))^2}{2\sigma_i^2}}$$

$$\log p(t|x, w, \beta) = \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma_i^2}} \right) - \sum_{i=1}^N \frac{1}{2\sigma_i^2} (t_i - y(x_i, w))^2$$

$$\log p(w|x, t, \beta, \alpha) = \log p(t|x, w, \beta) + \log p(w|\alpha) - \log p(t|x, \beta, \alpha)$$

does not depend on w

$$-\log p(w|x, t, \beta, \alpha) = -\log p(t|x, w, \beta) - \log p(w|\alpha) + \log p(t|x, \beta, \alpha)$$

$$w_{\text{MAP}} = \arg \min_w -\log p(w|x, t, \beta, \alpha)$$

$$= \arg \min_w \left[\frac{1}{2} \sum_{i=1}^N \frac{1}{\sigma_i^2} (t_i - y(x_i, w))^2 + \frac{\alpha}{2} w^T w \right]$$

c) from maximum likelihood, objective function is given by

$$\frac{1}{2} \sum_{i=1}^N \frac{1}{\sigma_i^2} (t_i - y(x_i, w))^2$$

$$\text{let } r_i = \frac{1}{\sigma_i^2} \quad (r_i > 0)$$

$$\Rightarrow E_D(w) = \frac{1}{2} \sum_{i=1}^N r_i (t_i - y(x_i, w))^2$$

↓
sum of squares error
function with weighting factor

let $y(x, w) = w^T \phi(x)$ be linear in terms of
basis functions $\phi_i(x)$

$$\Rightarrow E_D(w) = \frac{1}{2} \sum_{i=1}^N r_i (t_i - w^T \phi(x_i))^2$$

To minimize $E_D(w)$, $\frac{\partial}{\partial w} E_D(w) = 0$

$$\Rightarrow \frac{\partial}{\partial w} \left[\frac{1}{2} \sum_{i=1}^N r_i (t_i - w^T \phi(x_i))^2 \right] = 0$$

$$\frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial w} r_i (t_i - w^T \phi(x_i))^2 = 0$$

$$\frac{1}{2} \sum_{i=1}^N 2 r_i (t_i - w^T \phi(x_i)) \cdot \frac{\partial}{\partial w} (-w^T \phi(x_i)) = 0$$

$$\sum_{i=1}^N r_i (t_i - w^T \phi(x_i)) \frac{\partial}{\partial w} (-\phi(x_i)^T w) = 0$$

$$\sum_{i=1}^N r_i (t_i - w^T \phi(x_i)) \phi(x_i)^T = 0$$

$$\Rightarrow \sum_{i=1}^N \omega^T r_i \phi(x_i) \phi(x_i)^T = \sum_{i=1}^N r_i t_i \phi(x_i)^T$$

Taking transpose:

$$\left(\sum_{i=1}^N r_i \phi(x_i) \phi(x_i)^T \right) \omega = \sum_{i=1}^N r_i t_i \phi(x_i)$$

$$\text{Let } \Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \vdots & \vdots & & \vdots \\ \phi_0(x_N) & \dots & \dots & \phi_{M-1}(x_N) \end{bmatrix}_{N \times M}$$

~~Let R =~~

$\sum_{i=1}^N r_i \phi(x_i) \phi(x_i)^T$ can be expressed as

$$\begin{bmatrix} r_1 \phi_0(x_1) & r_2 \phi_0(x_2) & \dots & r_n \phi_0(x_n) \\ r_1 \phi_1(x_1) & r_2 \phi_1(x_2) & & \vdots \\ \vdots & \vdots & & \vdots \\ r_1 \phi_{M-1}(x_1) & r_2 \phi_{M-1}(x_2) & & r_n \phi_{M-1}(x_n) \end{bmatrix}_{M \times N} \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \vdots & \vdots & & \vdots \\ \phi_0(x_N) & \dots & \dots & \phi_{M-1}(x_N) \end{bmatrix}_{N \times M}$$

\downarrow

$$\Phi^T R$$

\downarrow

$$\Phi$$

$$\text{where } R = \begin{bmatrix} r_1 & 0 & 0 & \dots & 0 \\ 0 & r_2 & \dots & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \dots & r_n \end{bmatrix}_{N \times N} \rightarrow \text{diagonal matrix.}$$

$$\Rightarrow (\Phi^T R \Phi) \omega = \Phi^T R t$$

$$\Rightarrow \boxed{\omega = (\Phi^T R \Phi)^{-1} \Phi^T R t}$$

$$\text{where } t = \begin{bmatrix} t_1 \\ \vdots \\ t_n \end{bmatrix}_{N \times 1}$$