

3.3 Given the error function $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N r_n \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$
 To minimize it, we can take its derivative w.r.t \mathbf{w} and set it to 0

$$\begin{aligned} \frac{dE_D(\mathbf{w})}{d\mathbf{w}} &= \frac{d\left(\frac{1}{2} \sum_{n=1}^N r_n \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2\right)}{d\mathbf{w}} \\ &= \frac{1}{2} \sum_{n=1}^N r_n \frac{d\left(\{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2\right)}{d\mathbf{w}} \end{aligned}$$

Using 3.13 from PRML,

$$= \sum_{n=1}^N r_n \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T$$

Setting this gradient to zero gives:

$$\begin{aligned} 0 &= \sum_{n=1}^N r_n \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T \\ \implies 0 &= \sum_{n=1}^N (r_n t_n \phi(\mathbf{x}_n)^T - r_n \mathbf{w}^T \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T) \\ \implies 0 &= \sum_{n=1}^N r_n t_n \phi(\mathbf{x}_n)^T - \sum_{n=1}^N r_n \mathbf{w}^T \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \\ \implies 0 &= \sum_{n=1}^N r_n t_n \phi(\mathbf{x}_n)^T - \sum_{n=1}^N \mathbf{w}^T (\phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T) r_n \quad \text{since } r_n \text{ is a scalar} \\ \implies 0 &= \sum_{n=1}^N r_n t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \sum_{n=1}^N (\phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T) r_n \end{aligned}$$

Let \mathbf{R} be a diagonal matrix where the n^{th} diagonal value is given by r_n .

$$\begin{aligned} \implies 0 &= (\Phi^T \mathbf{R} \mathbf{t})^T - \mathbf{w}^T \Phi^T \mathbf{R} \Phi \\ \implies \mathbf{w}^T \Phi^T \mathbf{R} \Phi &= (\Phi^T \mathbf{R} \mathbf{t})^T \end{aligned}$$

Taking transpose of both sides, we get:

$$\implies \Phi^T \mathbf{R} \Phi \mathbf{w} = (\Phi^T \mathbf{R} \mathbf{t})$$

$$\implies \mathbf{w}_* = (\Phi^T \mathbf{R} \Phi)^{-1} (\Phi^T \mathbf{R} \mathbf{t})$$

Another method for solving the same would be:

Let $\sqrt{r_n} \phi(\mathbf{x}_n) = \phi'(\mathbf{x}_n)$, and $\sqrt{r_n} t_n = t'_n$.

Then, the above expression becomes

$$0 = \sum_{n=1}^N t'_n \phi'(\mathbf{x}_n)^T - \mathbf{w}^T \sum_{n=1}^N (\phi'(\mathbf{x}_n) \phi'(\mathbf{x}_n)^T)$$

Using 3.15 from PRML, we get a very similar result:

$$\mathbf{w}^* = (\Phi'^T \Phi')^{-1} \Phi'^T \mathbf{t}'$$

$$\text{where } \Phi' = \begin{bmatrix} \sqrt{r_1} \phi_0(\mathbf{x}_1) & \sqrt{r_1} \phi_1(\mathbf{x}_1) & \dots & \sqrt{r_1} \phi_{M-1}(\mathbf{x}_1) \\ \sqrt{r_2} \phi_0(\mathbf{x}_2) & \sqrt{r_2} \phi_1(\mathbf{x}_2) & \dots & \sqrt{r_2} \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \dots & \vdots \\ \sqrt{r_N} \phi_0(\mathbf{x}_N) & \sqrt{r_N} \phi_1(\mathbf{x}_N) & \dots & \sqrt{r_N} \phi_{M-1}(\mathbf{x}_N) \end{bmatrix}$$

and

$$\mathbf{t}' = \begin{bmatrix} \sqrt{r_1} t_1 \\ \sqrt{r_2} t_2 \\ \vdots \\ \sqrt{r_N} t_N \end{bmatrix}$$

Interpretations:

1. Data dependent noise variance : we assumed initially that the target variable t is given by $t_i = \mathbf{w}^T \phi(\mathbf{x}_i) + \epsilon$ for all i . In this case, this equation becomes $\sqrt{r_i} t_i = \sqrt{r_i} \mathbf{w}^T \phi(\mathbf{x}_i) + \epsilon \implies t_i = \mathbf{w}^T \phi(\mathbf{x}_i) + \frac{\epsilon}{\sqrt{r_i}}$. Since ϵ is a zero mean Gaussian random variable that represents noise, we can see that the noise is having data dependent variance where $\sigma'_i = (\frac{1}{\sqrt{r_i}})^2 \sigma = \frac{\sigma}{r_i}$.

2. Replicated data points : Data points with higher r_i value associated with them are replicated more in the dataset. This confirms the above result where higher r_i is associated with lower variance, as more data points that are replicated in the dataset imply reduced uncertainty around that data point.