

BLOCK II EXERCISES:

LINEAR REGRESSION

[SOLUTIONS]

Ex.2

$$i) L(x, y, w) = \frac{1}{2} \sum_{n=1}^N r_n (y_n - w^\top \phi(x_n))^2$$

* $L(x, y, w)$ can be written in matrix form

$$= \frac{1}{2} (y - \Phi w)^\top R (y - \Phi w)$$

* R is a diagonal matrix with $R_{nn} = r_n$

$$= \frac{1}{2} (y^\top - (\Phi w)^\top) R (y - \Phi w)$$

$$= \frac{1}{2} (y^\top - w^\top \Phi^\top) R (y - \Phi w)$$

$$= \frac{1}{2} (y^\top R y - w^\top \Phi^\top R y - y^\top R \Phi w + w^\top \Phi^\top R \Phi w)$$

* Notice that

$$[y^\top (R \Phi w)]^\top = (R \Phi w)^\top y = (\Phi w)^\top R^\top y = w^\top \Phi^\top R^\top y,$$

which happens to be a scalar

$$= \frac{1}{2} (y^\top R y - 2w^\top \Phi^\top R y + w^\top \Phi^\top R \Phi w)$$

Now we derivative above expression with respect to w

$$\begin{aligned} \frac{\partial L(w)}{\partial w} &= -\frac{\partial}{\partial w} (w^\top \Phi^\top R y) + \frac{1}{2} \frac{\partial}{\partial w} (w^\top \Phi^\top R \Phi w) \\ &= -\Phi^\top R y + \Phi^\top R \Phi w \end{aligned}$$

Finally

$$\frac{\partial L(w)}{\partial w} = 0 \rightarrow \Phi^\top R \Phi w = \Phi^\top R y \rightarrow \boxed{w = (\Phi^\top R \Phi)^{-1} \Phi^\top R y}$$

Matrix dimensions

w $\Delta \times 1$

$\phi(x_n)$ $\Delta \times 1$

Φ $N \times \Delta$

y $N \times 1$

R $N \times N$

ii) If we approach the problem as likelihood maximization under a Gaussian distribution, then R can be interpreted as a precision matrix, i.e. inverse of the covariance matrix.

$$\max \log N(y | \phi w, R) = \max -\frac{1}{2} (y - \phi w)^T R (y - \phi w)$$

so r_n can be interpreted as the precision parameter of the n -th sample (x_n, y_n) . ↑ inverse of the variance.

iii) r_n can represent the number of times (x_n, y_n) appears in our dataset.

Ex.3

We have $f_{\varepsilon}(x, w) = w_0 + \sum_{i=1}^D w_i(x_i + \varepsilon_i) = f(w, x) + \sum_{i=1}^D w_i \varepsilon_i$
 $\varepsilon_i \sim N(0, \sigma^2)$

We need to show

$$\min E_{\varepsilon}[L_{\varepsilon}(x, y, w)] = \min L(x, y, w) + \|w\|_2^2.$$

Let us define the loss with the noise

$$L_{\varepsilon}(x, y, w) = \frac{1}{2} \sum_{n=1}^N (\hat{y}_n - y_n)^2$$

$$= \frac{1}{2} \sum_n (\hat{y}_n^2 - 2\hat{y}_n y_n + y_n^2)$$

$$= \frac{1}{2} \sum_n \left(\underset{(1)}{\hat{y}_n^2} + \left(\underset{(2)}{\sum_i w_i \varepsilon_{ni}} \right)^2 + \underset{(3)}{2 \hat{y}_n \sum_i w_i \varepsilon_{ni}} \right. \\ \left. - \underset{(4)}{2 \hat{y}_n y_n} - \underset{(5)}{2 y_n \sum_i w_i \varepsilon_{ni}} + \underset{(6)}{y_n^2} \right)$$

If we take the expectation with respect to ε

(3) and (5) disappear $E[\varepsilon_{ni}] = 0$

$$E_{\varepsilon}[L_{\varepsilon}(x, y, w)] = \frac{1}{2} \sum_n (\hat{y}_n - y_n)^2 + \frac{1}{2} E_{\varepsilon} \left[\sum_n \left(\sum_i w_i \varepsilon_{ni} \right)^2 \right]$$

$\underbrace{\hspace{10em}}_{(1)}$ $\underbrace{\hspace{10em}}_{(2)}$

$$(2) \left(\sum_i w_i \varepsilon_{ni} \right)^2 = \sum_{i=1}^D w_i^2 \varepsilon_{ni}^2 + 2 \sum_{j=1}^D \sum_{k=1}^{j-1} w_j \varepsilon_{nj} w_k \varepsilon_{nk}$$

We know that $E[\varepsilon_i \varepsilon_j] = 0$ $i \neq j$, then

$$E_{\varepsilon} \left[\sum_n \left(\sum_i w_i \varepsilon_{ni} \right)^2 \right] = \sum_i w_i^2 \sigma^2$$

$$\sigma^2 = E[\varepsilon^2] + E[\varepsilon]^2$$

So finally

$$E_{\varepsilon}[L_{\varepsilon}(x, y, w)] = L(x, y, w) + \frac{1}{2} N \sum_i w_i^2 \sigma^2$$

(Ex. 5)

Recall
 LSE: $\hat{w} = (X^T X)^{-1} X^T Y \rightarrow \hat{\theta}(x) = x^T \hat{w} = \hat{\theta}$

Now observation
 $\text{Var}(\hat{\theta}) = E[(\hat{\theta} - E[\hat{\theta}])^2]$

$D = \{x_{n \times 1}, y_{n \times 1}\}$

$\text{Bias}(\hat{\theta}, \theta) = E[\hat{\theta}] - \theta$

i)

We need to proof $\text{Var}(\hat{\theta}) \leq \text{Var}(\tilde{\theta})$ where $\tilde{\theta} = c^T \hat{w} \neq \hat{\theta}$
 LSE \uparrow any linear unbiased estimator

Step 0 Assumptions

True model $y_{n \times 1} = \theta + \varepsilon = Xw + \varepsilon$, $\varepsilon \sim N(0, \sigma^2 I)$

$$\text{Var}(y) = \sigma^2 I$$

$$\text{Bias}(\hat{\theta}, \theta) = 0 \rightarrow E[\hat{\theta}] = \theta$$

cte training response.

Any other linear unbiased estimator $\sim \tilde{\theta} = c^T y$

Step 1 Is $\hat{\theta}$ unbiased? i.e. $E[\hat{\theta}] = \theta$

$$E[\hat{\theta}] = E[c^T y] = E[c^T (\theta + \varepsilon)]$$

$$= E[c^T \theta] + E[c^T \varepsilon] = c^T \theta = c^T Xw = c^T Xw = \alpha^T w = \theta$$

$\hat{\theta}$ is unbiased only if $c^T X = \alpha^T$

Step 2 $\text{Var}(\hat{\theta}) \leq \text{Var}(\tilde{\theta})$

Trick: $\tilde{\theta} = \hat{\theta} + (\hat{\theta} - \tilde{\theta})$

$$\text{Var}(\tilde{\theta}) = \dots$$

Hint: $\text{Var}(x_1 + x_2) = \sum_{i=1}^2 \sum_{j=1}^2 \text{cov}(x_i, x_j)$

Ex.6

$$i) p(w) = N(0, \tau I) \quad p(y|x, w) = N(x^T w, \sigma^2)$$

$$\hat{w} = \underset{w}{\operatorname{argmax}} p(w|x, y) = \underset{w}{\operatorname{argmax}} p(w)p(y|x, w)$$

$$= \underset{w}{\operatorname{argmax}} [\log p(w) + \log p(y|x, w)]$$

$$= \underset{w}{\operatorname{argmax}} \left[\frac{1}{2\sigma^2} \|w\|_2^2 + \frac{1}{2\tau^2} \|y - x^T w\|_2^2 \right]$$

$$\times (2\sigma^2) \quad = \underset{w}{\operatorname{argmax}} \left[\|w\|_2^2 + \frac{\sigma^2}{\tau^2} \|y - x^T w\|_2^2 \right]$$

$$ii) \lambda = \frac{\sigma^2}{\tau^2}$$