

1.

因為 x_n 和 x_m 是獨立的，所以有：

$$E[x_n x_m] = E[x_n] E[x_m]$$

由於 x_n 和 x_m 都來自於平均值為 μ 、方差為 σ^2 的高斯分布，因此有：

$$E[x_n] = E[x_m] = \mu$$

綜合上面兩個公式，得到：

$$E[x_n x_m] = E[x_n] E[x_m] = \mu \mu = \mu^2$$

如果 $n = m$ ，那麼：

$$E[x_n x_m] = E[x_n^2] = \text{Var}[x_n] + E[x_n]^2 = \sigma^2 + \mu^2$$

否則， $n \neq m$ ， x_n 和 x_m 獨立，所以它們的協方差為 0，即：

$$E[x_n x_m] = 0$$

所以有：

$$E[x_n x_m] = \mu^2 + I_{nm} \sigma^2 - (3)$$

當 $n = m$ 時， $I_{nm} = 1$ ，否則 $I_{nm} = 0$ ，得證。

為了證明 $E[\mu_{ML}] = \mu$ ，我們需要顯示最大近似估計平均值參數的期望值等於數據生成過程的真實平均值。

對於具有未知平均值 μ 和已知方差 σ^2 的高斯分布，其概似函數可以寫為：

$$L(\mu | x_1, x_2, \dots, x_n) = (2\pi\sigma^2)^{-\left(\frac{n}{2}\right)} e^{-(1/2\sigma^2)\sum(x_i - \mu)^2}$$

對數近似函數為：

$$\begin{aligned} \ln L(\mu | x_1, x_2, \dots, x_n) &= (2\pi\sigma^2)^{-\left(\frac{n}{2}\right)} e^{-\left(\frac{1}{2\sigma^2}\right)\sum(x_i - \mu)^2} \\ &= -\left(\frac{n}{2}\right) \ln(2\pi) - \left(\frac{n}{2}\right) \ln(\sigma^2) - \left(\frac{1}{2\sigma^2}\right) \sum(x_i - \mu)^2 \end{aligned}$$

為了找到 μ 的最大概似估計，我們對對數概似函數對 μ 求導數並將其設置為零：

$$\frac{d}{du} \ln L(\mu | x_1, x_2, \dots, x_n) = \left(\frac{1}{\sigma^2}\right) \sum(x_i - \mu) = 0$$

解出 μ ，我們得到：

$$\mu_{ML} = \frac{1}{n} \sum x_i$$

μ 的最大概似估計值期望可以通過對 μ_{ML} 取期望值得到：

$$E[\mu_{ML}] = E\left[\frac{1}{n} \sum x_i\right] = \frac{1}{n} \sum E[x_i] = \mu$$

在這裡，我們使用了每個數據點 x_i 的期望值等於真實平均值 μ 的事實。

為了證明 $E[\sigma_{ML}^2] = \frac{N-1}{N} \sigma^2$ ，可以通過將 μ 的最大概似估計值代入對數概似函數

並對 σ^2 進行最大化來找到變量參數的最大概似估計：

$$\ln L(\sigma^2 | x_1, x_2, \dots, x_n) = -\left(\frac{n}{2}\right) \ln(2\pi) - \left(\frac{n}{2}\right) \ln(\sigma^2) - \left(\frac{1}{2\sigma^2}\right) \sum (x_i - \mu)^2$$

$$\frac{d}{d\sigma^2} \ln L(\sigma^2 | x_1, x_2, \dots, x_n) = -\left(\frac{n}{2}\right) \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum (x_i - \mu)^2 = 0$$

解出 σ^2 ，我們得到：

$$\sigma_{ML}^2 = \frac{1}{n} \sum (x_i - \mu)^2$$

σ^2 的最大概似估計期望可以通過對 σ_{ML}^2 取期望值得到：

$$E[\sigma_{ML}^2] = E\left[\frac{1}{n} \sum (x_i - \mu)^2\right]$$

使用 $\text{Var}(x_i)$ ，因此，我們得出 MLE 的期望值是偏小估計量，即

$$E[\sigma_{ML}^2] = \frac{N-1}{N} \sigma^2$$

得證。

2.

我們可以使用期望的線性來證明兩個獨立隨機向量 \mathbf{a} 和 \mathbf{b} 的和 $\mathbf{y} = \mathbf{a} + \mathbf{b}$ 的平均值由每個變量的平均值分別組成，

即：

$$E[\mathbf{y}] = E[\mathbf{a} + \mathbf{b}] = E[\mathbf{a}] + E[\mathbf{b}]$$

其中 $E[\mathbf{y}]$ 是隨機變量 \mathbf{y} 的期望值， $E[\mathbf{a}]$ 和 $E[\mathbf{b}]$ 分別是隨機變量 \mathbf{a} 和 \mathbf{b} 的期望值，這是由於隨機變量的和之期望值等於它們各自期望值的總和。

為了顯示 \mathbf{y} 的協方差矩陣，由 \mathbf{a} 和 \mathbf{b} 的協方差矩陣之和給出，我們可以使用協方差的定義，即：

$$\text{cov}(\mathbf{y}) = E[(\mathbf{y} - E[\mathbf{y}])(\mathbf{y} - E[\mathbf{y}])^T]$$

其中 $\text{cov}(\mathbf{y})$ 是隨機向量 \mathbf{y} 的協方差矩陣， $E[\mathbf{y}]$ 是 \mathbf{y} 的平均值。

展開上式，並使用期望的線性性質，我們得到：

$$\begin{aligned} \text{cov}(\mathbf{y}) &= E[(\mathbf{a} - E[\mathbf{a}] + \mathbf{b} - E[\mathbf{b}])(\mathbf{a} - E[\mathbf{a}] + \mathbf{b} - E[\mathbf{b}])^T] = \\ &E[(\mathbf{a} - E[\mathbf{a}])(\mathbf{a} - E[\mathbf{a}])^T] + E[(\mathbf{a} - E[\mathbf{a}])(\mathbf{b} - E[\mathbf{b}])^T] \\ &+ E[(\mathbf{b} - E[\mathbf{b}])(\mathbf{a} - E[\mathbf{a}])^T] + E[(\mathbf{b} - E[\mathbf{b}])(\mathbf{b} - E[\mathbf{b}])^T] \\ &= \text{cov}(\mathbf{a}) + \text{cov}(\mathbf{b}) \end{aligned}$$

其中 $\text{cov}(\mathbf{a})$ 和 $\text{cov}(\mathbf{b})$ 分別是 \mathbf{a} 和 \mathbf{b} 的協方差矩陣。這是由於兩個獨立隨機變量的和的協方差等於它們各自的協方差之和。因此 \mathbf{y} 的協方差矩陣，由 \mathbf{a} 和 \mathbf{b} 的協方差矩陣之和給出。

3.

要證明 $\sigma^2_{N+1}(x) \leq \sigma^2_N(x)$ ，需要證明

$$\phi(x)^T S_{N+1} \phi(x) \leq \phi(x)^T S_N \phi(x)$$

我們使用 Woodbury matrix identity from Appendix C in [Bishop, 2006]

$$S_{N+1} = [S_N^{-1} + \beta \phi(x) \phi(x)^T]^{-1} = S_N - \frac{S_N \phi(x) \phi(x)^T S_N}{1 + \phi(x)^T S_N \phi(x)}$$

將這方程式帶入上面不等式，可以得到

$$\phi(x)^T \frac{S_N \phi(x) \phi(x)^T S_N}{1 + \phi(x)^T S_N \phi(x)} \phi(x) \geq 0$$

只要協方差矩陣 S_N 是半正定， $\gamma = \phi(x)^T S_N \phi(x) \geq 0$ ，上述不等式會變成

$$\frac{\gamma^2}{1 + \gamma} \geq 0$$

隨著數據集大小的增加，線性回歸函數相關的不確定性會減少。

4.

這個問題提到獨立的噪音被添加到每個維度的輸入變量 x_n ，因此我們的新模型變成，

$$\begin{aligned} y'(x_n, \omega) &= \omega_0 + \sum_{d=1}^D \omega_d (x_{nd} + \epsilon_{nd}) \\ &= \omega_0 + \sum_{d=1}^D \omega_d x_{nd} + \sum_{d=1}^D \omega_d \epsilon_{nd} = y(x_n, \omega) + \sum_{d=1}^D \omega_d \epsilon_{nd} \end{aligned}$$

其中 the noise ϵ_{nd} is independent across both the n and d indices，所以新的 error function is

$$\begin{aligned} E'_D(\omega) &= \frac{1}{2} \sum_{n=1}^N \{y'(x_n, \omega) - t_n\}^2 \\ &= \frac{1}{2} \sum_{n=1}^N \{y(x_n, \omega) + \sum_{d=1}^D \omega_d \epsilon_{nd} - t_n\}^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left\{ (y(x_n, \omega) - t_n)^2 + 2(y(x_n, \omega) - t_n) \left(\sum_{d=1}^D \omega_d \epsilon_{nd} \right) + \left(\sum_{d=1}^D \omega_d \epsilon_{nd} \right)^2 \right\} \end{aligned}$$

取期望值，並使用期望的線性性質，得到

$$\begin{aligned} \mathbb{E}[\mathbf{E}'_D(\boldsymbol{\omega})] &= \frac{1}{2} \sum_{n=1}^N \left\{ (\mathbf{y}(x_n, \boldsymbol{\omega}) - t_n)^2 + 2(\mathbf{y}(x_n, \boldsymbol{\omega}) - t_n) \left(\sum_{d=1}^D \omega_d \mathbb{E}[\epsilon_{nd}] \right) \right. \\ &\quad \left. + \mathbb{E} \left[\left(\sum_{d=1}^D \omega_d \epsilon_{nd} \right)^2 \right] \right\} \end{aligned}$$

$\mathbb{E}[\epsilon_{nd}]$ is 0, so the second term disappears. Now we look at the third term

$$\begin{aligned} \mathbb{E} \left[\left(\sum_{n=1}^N \omega_d \epsilon_{nd} \right)^2 \right] &= \mathbb{E} \left[\sum_{d=1}^D \sum_{d'=1}^D \omega_d \omega_{d'} \epsilon_{nd} \epsilon_{nd'} \right] \\ &= \sum_{d=1}^D \sum_{d'=1}^D \omega_d \omega_{d'} \mathbb{E}[\epsilon_{nd} \epsilon_{nd'}] \\ &= \sum_{d=1}^D \sum_{d'=1}^D \omega_d \omega_{d'} \delta_{dd'} = \sum_{d=1}^D \omega_d^2 \end{aligned}$$

Using these results, we get

$$\mathbb{E}[\mathbf{E}'_D(\boldsymbol{\omega})] = \frac{1}{2} \sum_{n=1}^N \left\{ (\mathbf{y}(x_n, \boldsymbol{\omega}) - t_n)^2 + \sum_{d=1}^D \omega_d^2 \right\} = \mathbf{E}_D(\boldsymbol{\omega}) + \frac{N}{2} \sum_{d=1}^D \omega_d^2$$

and we see that we get a L_2 regularization term without the bias parameter ω_0 , as desired.