(1)  As we know : $\Gamma(2,1)$  PPF  :  $P(z) = z \exp(-z)$ , $z > 0$ .

$$P_x(x) = \frac{\lambda^x}{x!} \exp(-\lambda) \quad , \lambda > 0 \ , \ x = 0, 1, 2 \dots$$

because $x_1, x_2 \dots x_d$ is an iid

MAP :  $\underset{\lambda}{\arg} \ P(w|x) = \underset{\lambda}{\arg \max} \ \dfrac{P(x|z) \ P(z)}{P(x)}$ .

$$\propto \underset{\lambda}{\arg \max} \ P(x|z) \ P(z) .$$

$$= \underset{\lambda}{\arg \max} \ \prod_{i=1}^{h} \left( \frac{\lambda^{x_i}}{x_i!} \exp(-\lambda) \right) \cdot z \exp(-z)$$

$$= \underset{\lambda}{\arg \max} \ \sum_{i=1}^{n} [x_i \ln\lambda - \lambda - \ln x_i! + \ln z - z$$

$$= \underset{\lambda}{\arg \max} \ [ \ln\lambda \sum_{i=1}^{h} x_i - n\lambda - \sum_{i=1}^{n} \ln(x_i!) + \ln z - z ] \rightarrow ①$$

$$\frac{\partial \theta}{\partial \lambda} = \frac{\sum_{i=1}^{h} x_i}{\lambda} - n = 0 \quad \Rightarrow \quad \sum_{i=1}^{h} x_i = n\lambda . \ \rightarrow \lambda = \frac{1}{n} \sum_{i=1}^{h} x_i > 0.$$

$$\frac{\partial^2 \theta}{\partial \lambda^2} = \frac{-\sum_{i=1}^{h} x_i}{\lambda^2} < 0 .$$

$\therefore \ \lambda = \frac{1}{n} \sum_{i=1}^{h} x_i$ . is a value that maximizes the posterior.

2. $\quad P_{X|\lambda}(X|\lambda) = \frac{\lambda^X}{X!}\exp(-\lambda) \qquad X=0, 1, 2 \cdots$

$$L(x|\theta) = \prod_{i=1}^{n} P(X_i|\lambda) = \sum_{i=1}^{h} \ln P(X_i|\lambda)$$

$$\hat{\lambda} = \arg\max_{\lambda} \sum_{i=1}^{n} [X_i \ln\lambda - \lambda - \ln(x_i!)]$$

$\quad$ set $\mu = \frac{1}{h}\sum_{i=1}^{h} X_i$

$$= \arg\max_{\lambda} [\ln(\lambda)\cdot n\mu - n\lambda - \sum_{i=1}^{n}\ln(x_i!)] \qquad \cdots \quad \textcircled{1}$$

$$\frac{\partial\theta}{\partial\lambda} = \frac{n\mu}{\lambda} - n \qquad = 0 \qquad \Rightarrow \quad \lambda = \mu.$$

$$\frac{\partial^2\theta}{\partial\lambda^2} = -\frac{n\mu}{\lambda^2} = -\frac{n}{\mu} < 0 \cdot \qquad \leftarrow \quad \text{maxium}.$$

$$\Rightarrow \quad \hat{\lambda} = \frac{1}{h}\sum_{i=1}^{h} X_i.$$

CLT: given enough samples, data will follow normal distribution.

$N(\mu, 6^2)$. which means

$$\hat{\lambda} = \arg\max_{\lambda} N(\mu, 6^2) = \mu = \frac{1}{h}\sum_{i=1}^{N} X_i.$$

$\quad$ From CLT, $\quad \hat{\lambda} = \frac{1}{h}\sum_{i=1}^{n} X_i$

(3)

$$Y = \beta_0 + \beta_1 x + \cdots \beta_p x_p + \varepsilon \quad , \quad \varepsilon \sim N(0, \sigma^2)$$

$$\Rightarrow Y = x\beta + \varepsilon \quad , \quad \text{where} \quad x = [1, x_1, \cdots x_p]$$

As we know, least squares: $\beta = (x^T x)^{-1} x^T y$.

$$\text{MLE}: \quad \hat{\beta} = \underset{\beta}{\arg\max} \; P(Y|D) \qquad D = \{x_1 \cdots x_p\}.$$

$$= \underset{\beta}{\arg\max} \; \prod_{i=1}^{n} P(Y|x_i)$$

$$\varepsilon = y - x\beta = N(0, \sigma^2)$$

$$\hat{\beta} = \underset{\beta}{\arg\max} \; \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi \sigma^2}} \exp\left(-\frac{1}{2}\left(\frac{y_i - x\beta_i}{\sigma}\right)^2\right).$$

$$= \underset{\beta}{\arg\max} \; \sum_{i=1}^{n} \ln \frac{1}{\sqrt{2\pi \sigma^2}} - \frac{1}{2}\left(\frac{y_i - x\beta_i}{\sigma}\right)^2 \quad \cdots \cdots \; \textcircled{1}$$

$$\frac{\partial \textcircled{1}}{\partial \beta} = \sum_{i=1}^{n} \frac{x^T}{\sigma}\left(\frac{y_i - x\beta_i}{\sigma}\right) = 0 \; \Rightarrow \; x^T x \beta = x^T y \; .$$

$$\therefore \hat{\beta} = (x^T x)^{-1} x^T y$$

$$\frac{\partial^2 \textcircled{1}}{\partial \beta^2} = \frac{-x^T x}{\sigma^2} < 0 \; , \; \leftarrow \text{maximum.}$$

(4) $Y = \beta_0 + \sum_{i=1}^{P} \beta_i x_i + \varepsilon$ . $\varepsilon \sim N(0, \delta^2)$.

$Y = x\hat{\beta} + \varepsilon$      $x = [1, x_1, \cdots x_p]$

$\varepsilon = y - x\hat{\beta} \sim N(0, \delta^2)$

$P_{Y/x}(Y|x_i) = \frac{1}{\sqrt{2\pi}\delta^2} \exp\left(-\frac{1}{2}\left(\frac{y - x_i\beta_i}{\delta}\right)^2\right)$.

MAP: $\hat{\beta} = \underset{\beta}{\arg\max}\left[\prod_{i=1}^{P} \frac{1}{\sqrt{2\pi}\delta} \exp\left(-\frac{1}{2}\left(\frac{y - x_i\beta_i}{\delta}\right)^2\right)\right] \cdot \sqrt{\frac{\lambda}{2\pi\delta^2}} \exp\left(-\frac{\lambda}{2}\left(\frac{\beta_i}{\delta}\right)^2\right)$

$\Rightarrow \hat{\beta} = \underset{\beta}{\arg\max} \sum_{i=1}^{P}\left[-\frac{1}{2}\ln(2\pi) - \ln\delta - \frac{1}{2}\left(\frac{y - x_i\beta_i}{\delta}\right)^2\right] + \frac{1}{2}\ln\lambda - \frac{1}{2}\ln(2\pi) - \ln\delta - \frac{\lambda}{2}\left(\frac{\beta_i}{\delta}\right)^2$

$= \arg\max\left[\frac{P+1}{2}\ln 2\pi + (P+1)\ln\delta - \frac{1}{2}\ln(\lambda) + \frac{\lambda}{2}\left(\frac{\beta_i}{\delta}\right)^2 + \frac{1}{2}\left(\frac{y - x_i\beta_i}{\delta}\right)^2\right]$ $\cdots$ ①

$\frac{\partial \Theta}{\partial \beta} = \frac{\lambda\beta_i}{\delta^2} - \frac{x^T(y - x\beta)}{\delta^2} = 0 \Rightarrow (x^Tx + \lambda)\beta = x^Ty$ .

$\Rightarrow \hat{\beta} = (x^Tx + \lambda)^{-1} x^Ty$.

Based on probabalistic structure, $\lambda$ should add bias to estimate.

$\Rightarrow \hat{\beta} = (x^Tx + \lambda\Sigma)^{-1} x^Ty$.

(I) $Y = X\beta + \varepsilon$, $\quad \varepsilon \sim N(0, \sigma^2)$.

$P(\beta_i) = lap(0, \frac{\sigma^2}{\lambda}) = \frac{\lambda}{2\sigma^2} \exp\left(\frac{-|\beta_i|}{\sigma^2}\right)$

$\varepsilon = y - X\beta \sim N(0, \sigma^2)$.

$P_{rix}(Y|X_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\left(\frac{(y-X_i\beta_i)^2}{\sigma}\right)\right)$

MAP: $\hat{\beta} = \arg\max_\beta \left[\prod_{i=1}^{P} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2}\left(\frac{y-X_i\beta_i}{\sigma}\right)^2\right)\right] \cdot \frac{\lambda}{2\sigma^2} \exp\left(-\frac{|\beta_i|}{\sigma^2}\right)$

$\hat{\beta} = \arg\min_\beta \sum_{i=1}^{P}\left[\frac{1}{2}\ln(2\pi) + \ln(\sigma) + \frac{1}{2}\left(\frac{y-X_i\beta_i}{\sigma}\right)^2 + \ln\left(\frac{\lambda}{2\sigma^2}\right) + \frac{|\beta_i|}{\sigma}\right]$

$\nabla\beta = \begin{cases} \frac{1}{\sigma^2}\left[-X^T(Y-X\beta) + 1\right] & \beta_i \geq 0. \\ \frac{1}{\sigma^2}\left[-X^T(Y-X\beta) - 1\right] & \text{otherwise} \end{cases} = 0$

$\beta_i \geq 0:$ $(X^TX)\beta = X^TY - \frac{1}{\sigma^2}$ $\Rightarrow$ $\beta = (X^TX)^{-1}(X^TY - \frac{1}{\sigma^2})$

$\beta_i < 0:$ $(X^TX)\beta = X^TY + \frac{1}{\sigma^2}$ $\Rightarrow$ $\beta = (X^TX)^{-1}(X^TY + \frac{1}{\sigma^2})$

$\nabla_\beta^2 \beta = X^TX > 0$ $\Rightarrow$ positive.

$\Rightarrow$ therefore it's maximum.

$\hat{\beta} = \begin{cases} (X^TX)^{-1}(X^TY - \frac{1}{\sigma^2}) & \beta_i > 0. \\ (X^TX)^{-1}(X^TY + \frac{1}{\sigma^2}) & \text{otherwise.} \end{cases}$

(b) As we know $X = U\Sigma V^T$

a. $\hat{\beta}^{Ridge} = \underset{\beta}{\text{argmin}} \; RSS + \lambda \sum_{j=1}^{P} \beta_j^2 = \underset{\beta}{\text{argmin}} \; \|y - X\beta\|^2 + \lambda \|\beta\|^2$

$= \underset{\beta}{\text{argmin}} \; (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$

$= \underset{\beta}{\text{argmin}} \; y^T y - (X\beta)^T y + y^T X\beta + \beta^T (X^T X + \lambda I)\beta \quad \leftarrow \; ①$

$\frac{\partial \theta}{\partial \beta} = -X^T y - (y^T X)^T + 2\beta (X^T X + \lambda I) = 0$

$\Rightarrow \quad \beta (X^T X + \lambda I) = X^T y$

$\Rightarrow \quad \hat{\beta}^{Ridge} = (X^T X + \lambda I)^{-1} X^T y$

$\Rightarrow \hat{\beta}^{Ridge} = ((U\Sigma V^T)^T (U\Sigma V^T) + \lambda I)^{-1} (V\Sigma^T U^T) y$

$= (V\Sigma^T U^T U \Sigma V^T + \lambda I)^{-1} (V\Sigma^T U^T) y$

$\because \quad \Sigma^T = \Sigma$

$\Rightarrow \hat{\beta}^{Ridge} = (\Sigma^T \Sigma + \lambda I)^{-1} (V\Sigma V^T) y$

$\hat{y} - X\hat{\beta}^{Ridge} = U\Sigma V^T (\Sigma^T \Sigma + \lambda I)^{-1} (V\Sigma U^T) y$

$= \sum_{j=1}^{P} \sigma_j u_j v_j^T (\sigma_j^2 + \lambda)^{-1} v_j \sigma_j u_j^T y$

$= \sum_{j=1}^{P} \frac{\sigma_j^2}{\sigma_j^2 + \lambda} u_j v_j^T v_j u_j^T y$

$= \sum_{j=1}^{P} u_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} u_j^T y$

(b)

$$X(X^TX + \lambda I)^{-1}X^T = V\Sigma U^T (V\Sigma^T U^T U\Sigma V^T + \lambda I)^{-1} V\Sigma U^T$$

$$= V\Sigma U^T (V\Sigma^T \Sigma V^T + \lambda I)^{-1} V^T\Sigma U$$

$$= V\Sigma U^T (\Sigma^T\Sigma (VV^T) + \lambda I)^{-1} V\Sigma U^T$$

$$= \sum_{i=1}^{P} U_i \frac{\sigma_i^2}{\sigma_i^2 + \lambda} U_i^T$$

$$\Rightarrow \text{tr}\left(\sum_{i=1}^{P} U_i \left(\frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right) U_i^T\right) = \sum_{j=1}^{n}\sum_{i=1}^{P} \left(U_i \frac{\sigma_i^2}{\sigma_i^2+\lambda} U_i^T\right)_{jj} \cdot$$

$$= \sum_{j=1}^{n}\sum_{i=1}^{P} \left(\frac{\sigma_i^2}{\sigma_i^2+\lambda} U_i U_i^T\right)_{jj}$$

$$= \sum_{j=1}^{n} \frac{\sigma_j^2}{\sigma_j^2 + \lambda}$$

$$\therefore \text{tr}\left(X(X^TX + \lambda I)^{-1} X^T\right) = \sum_{j=1}^{n} \frac{\sigma_j^2}{\sigma_j^2 + \lambda}$$