



# Tarea 3

## Ciencia de Datos.

Alfredo Bistrain  
Jesús Salazar

October 14, 2025

### Ejercicio

#### 1 Regresión lineal ordinaria (OLS)

- Derivación del estimador OLS: Partiendo del modelo clásico

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

demuestre que el estimador de Mínimos Cuadrados Ordinarios es

$$\hat{\beta} = (X^\top X)^{-1} X^\top y$$

siempre que  $X^\top X$  sea invertible.

- Propiedades del estimador: Calcule explícitamente:

$$E[\hat{\beta}], \quad \text{Var}(\hat{\beta}).$$

Concluya que  $\hat{\beta}$  es insesgado y eficiente dentro de la clase de estimadores lineales (teorema de Gauss-Markov).

#### Solución:

Considere el modelo de regresión lineal,

$$\mathbf{Y} = X\beta + \epsilon$$

donde  $\epsilon \sim N(\mathbf{0}, \sigma^2 I_n)$ . Se busca al valor de  $\hat{\beta}$  que minimice la suma de los errores al cuadrado, es decir,

$$\begin{aligned} SS_{\text{Res}}(\beta) &= \sum_{i=1}^n (y_i - [X\beta]_i)^2 \\ &= (\mathbf{y} - X\beta)^\top (\mathbf{y} - X\beta) \\ &= \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top X\beta - \beta^\top X^\top \mathbf{y} + \beta^\top X^\top X\beta \\ &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top X\beta + \beta^\top X^\top X\beta. \end{aligned}$$

Derivando la expresión anterior con respecto a  $\beta$  y notando que  $t^\top X$  es una matriz

simétrica, se tiene,

$$\begin{aligned}\frac{\partial SS_{\text{Res}}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= (-2\mathbf{y}^\top X)^\top + 2X^\top X\boldsymbol{\beta} \\ &= -2X^\top \mathbf{y} + 2X^\top X\boldsymbol{\beta}.\end{aligned}$$

Igualando a cero la expresión anterior, se puede observar que el estimador que se busca  $\bar{\boldsymbol{\beta}}$  debe satisfacer

$$X^\top X\bar{\boldsymbol{\beta}} = X^\top \mathbf{y}.$$

Este sistema de ecuaciones son las ecuaciones normales para el modelo de regresión lineal. Para poder despejar  $\bar{\boldsymbol{\beta}}$  es necesario que  $X^\top X$  sea invertible. Esto es que tenga rango completo. En este caso,

$$\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{y}.$$

Sabemos que  $\hat{\boldsymbol{\beta}}$  minimiza  $SS_{\text{Res}}(\boldsymbol{\beta})$  si y solo si la matriz Hessiana de  $SS_{\text{Res}}$  es positiva definida.

$$\begin{aligned}H &= \frac{\partial^2 SS_{\text{Res}}}{\partial \boldsymbol{\beta}^2} \\ &= (2X^\top X)^\top \\ &= 2X^\top X.\end{aligned}$$

que es positiva definida si  $\text{rango}(X) = p$ .

(2)

Para solucionar este inciso, bastan observar que Como  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 I_n)$ , entonces  $\mathbf{Y} \sim N(X\boldsymbol{\beta}, \sigma^2 I_n)$ , esto por propiedades de la distribución normal multivariada, por otro lado, observe que,

$$\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbb{E}\left((X^\top X)^{-1} X^\top \mathbf{Y}\right) = (X^\top X)^{-1} X^\top \mathbb{E}(\mathbf{Y}) = (X^\top X)^{-1} X^\top X\boldsymbol{\beta} = \boldsymbol{\beta}.$$

Del mismo modo, la varianza de  $\hat{\boldsymbol{\beta}}$  se puede obtener como,

$$\begin{aligned}\text{Var}(\hat{\boldsymbol{\beta}}) &= \text{Var}\left[(X^\top X)^{-1} X^\top \mathbf{y}\right] \\ &= (X^\top X)^{-1} X^\top \text{Var}(\mathbf{y}) \left[(X^\top X)^{-1} X^\top\right]^\top \\ &= (X^\top X)^{-1} X^\top X (X^\top X)^{-1} \sigma^2 = (X^\top X)^{-1} \sigma^2.\end{aligned}$$

De este modo,

$$\hat{\boldsymbol{\beta}} \sim N\left(\boldsymbol{\beta}, (X^\top X)^{-1} \sigma^2\right),$$

y por el teorema de Gauss-Markov, se tiene que  $\hat{\boldsymbol{\beta}}$  es un estimador insesgado y eficiente.

**Ejercicio****2 Regresión lineal bayesiana (prior conjugado)**

1. Prior conjugado: Suponga un prior conjugado:

$$\beta \mid \sigma^2 \sim \mathcal{N}(\beta_0, \sigma^2 V_0), \quad \sigma^2 \sim \text{Inv-Gamma}(a_0, b_0).$$

2. Distribución posterior: Derive los parámetros posteriores ( $\beta_n, V_n, a_n, b_n$ ) y escriba la forma explícita de la posterior conjunta

$$p(\beta, \sigma^2 \mid y)$$

3. Distribuciones marginales: Identifique las distribuciones marginales de  $\beta$  y de  $\sigma^2$ .

**Solución:**

(1)

Considere el modelo de regresión lineal,

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

donde  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 I_n)$ . Además, suponga que los parámetros son variables aleatorias con distribuciones a priori,

$$\begin{aligned} \beta \mid \sigma^2 &\sim \mathcal{N}(\beta_0, \sigma^2 V_0) \\ \sigma^2 &\sim \text{Inv-Gamma}(a_0, b_0), \end{aligned}$$

donde  $V_0, \beta_0, a_0, b_0$  son hiperparámetros de las distribuciones. De este modo, la función de verosimilitud está dada por,

$$p(y \mid \beta, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} (y - \mathbf{X}\beta)^\top (y - \mathbf{X}\beta) \right\}.$$

Observe que la distribución a priori conjunta de los parámetros se puede obtener como,

$$\begin{aligned} p(\beta, \sigma^2) &= p(\beta \mid \sigma^2)p(\sigma^2) \\ &= (2\pi\sigma^2)^{-k/2} |V_0|^{-1/2} \exp \left\{ -\frac{1}{2\sigma^2} (\beta - \beta_0)^\top V_0^{-1} (\beta - \beta_0) \right\} \\ &\quad \times \frac{b_0^{a_0}}{\Gamma(a_0)} (\sigma^2)^{-a_0-1} \exp \left( -\frac{b_0}{\sigma^2} \right). \end{aligned}$$

Por otro lado, se sabe que la distribución posterior es proporcional a la expresión,

$$\begin{aligned}
p(\beta, \sigma^2 | y) &\propto p(y | \beta, \sigma^2)p(\beta | \sigma^2)p(\sigma^2) \\
&\propto (\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} (y - X\beta)^\top (y - X\beta) \right\} \\
&\quad \times (\sigma^2)^{-k/2} \exp \left\{ -\frac{1}{2\sigma^2} (\beta - \beta_0)^\top V_0^{-1} (\beta - \beta_0) \right\} \\
&\quad \times (\sigma^2)^{-a_0-1} \exp \left( -\frac{b_0}{\sigma^2} \right).
\end{aligned}$$

Note que si se juntan las dos primeras exponenciales y se factoriza el término  $-\frac{1}{2\sigma^2}$ , los términos dentro de la exponencial se pueden reescribir como,

$$\begin{aligned}
&(y - X\beta)^\top (y - X\beta) + (\beta - \beta_0)^\top V_0^{-1} (\beta - \beta_0) \\
&= y^\top y - 2\beta^\top X^\top y + \beta^\top X^\top X\beta + \beta^\top V_0^{-1}\beta - 2\beta^\top V_0^{-1}\beta_0 + \beta_0^\top V_0^{-1}\beta_0 \\
&= \beta^\top (X^\top X + V_0^{-1})\beta - 2\beta^\top (X^\top y + V_0^{-1}\beta_0) + y^\top y + \beta_0^\top V_0^{-1}\beta_0.
\end{aligned}$$

Para simplificar la notación, defina,

$$\begin{aligned}
A &= X^\top X + V_0^{-1} \\
b &= X^\top y + V_0^{-1}\beta_0 \\
c &= y^\top y + \beta_0^\top V_0^{-1}\beta_0
\end{aligned}$$

Observe que la última expresión se puede reescribir como,

$$\begin{aligned}
&\beta^\top A\beta - 2\beta^\top b + c \\
&= \beta^\top A\beta - 2\beta^\top b + b^\top A^{-1}b - b^\top A^{-1}b + c \\
&= (\beta - A^{-1}b)^\top A(\beta - A^{-1}b) - b^\top A^{-1}b + c \\
&= (\beta - A^{-1}b)^\top A(\beta - A^{-1}b) + (c - b^\top A^{-1}b).
\end{aligned}$$

Regresando esta expresión a ecuación original se tiene,

$$\begin{aligned}
p(\beta, \sigma^2 | y) &\propto (\sigma^2)^{-(n+k)/2-a_0-1} \exp \left\{ -\frac{1}{2\sigma^2} [(\beta - A^{-1}b)^\top A(\beta - A^{-1}b) + (c - b^\top A^{-1}b)] - \frac{b_0}{\sigma^2} \right\} \\
&= (\sigma^2)^{-k/2} \exp \left\{ -\frac{1}{2\sigma^2} (\beta - A^{-1}b)^\top A(\beta - A^{-1}b) \right\} \\
&\quad \times (\sigma^2)^{-n/2-a_0-1} \exp \left\{ -\frac{1}{\sigma^2} \left[ b_0 + \frac{1}{2}(c - b^\top A^{-1}b) \right] \right\}.
\end{aligned}$$

Así, multiplicando con las constantes adecuadas para completar las distribuciones correspondientes se tiene,

$$\begin{aligned}
p(\beta, \sigma^2 | y) &\propto \frac{1}{(2\pi\sigma^2)^{k/2}|A^{-1}|^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\beta - A^{-1}b)^\top A(\beta - A^{-1}b) \right\} \\
&\quad \times \frac{\left[ b_0 + \frac{1}{2}(c - b^\top A^{-1}b) \right]^{a_0+n/2}}{\Gamma(a_0 + n/2)} (\sigma^2)^{-(a_0+n/2)-1} \exp \left\{ -\frac{1}{\sigma^2} \left[ b_0 + \frac{1}{2}(c - b^\top A^{-1}b) \right] \right\}.
\end{aligned}$$

...

De esta manera, si se define a

$$\begin{aligned} V_n &= A^{-1} = (X^\top X + V_0^{-1})^{-1} \\ \beta_n &= A^{-1}b = V_n(X^\top y + V_0^{-1}\beta_0) \\ a_n &= a_0 + \frac{n}{2} \\ b_n &= b_0 + \frac{1}{2} [c - \beta_n^\top V_n^{-1}\beta_n] \\ &= b_0 + \frac{1}{2} [y^\top y + \beta_0^\top V_0^{-1}\beta_0 - \beta_n^\top V_n^{-1}\beta_n], \end{aligned}$$

se tiene que la distribución a posterior es una distribución normal-gamma inversa con parámetros,  $(\beta_n, \sigma^2 V_n, a_n, b_n)$

$$\begin{aligned} p(\beta, \sigma^2 | y) &= \frac{1}{(2\pi\sigma^2)^{k/2}|V_n|^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(\beta - \beta_n)^\top V_n^{-1}(\beta - \beta_n) \right\} \\ &\quad \times \frac{b_n^{a_n}}{\Gamma(a_n)} (\sigma^2)^{-a_n-1} \exp \left( -\frac{b_n}{\sigma^2} \right). \end{aligned}$$

Por último, observe que la distribución a posterior marginal de parámetro  $\sigma^2$  se puede encontrar a partir de la definición,

$$\begin{aligned} p(\sigma^2 | y) &= \int p(\beta, \sigma^2 | y) d\beta \\ &= \int \frac{1}{(2\pi\sigma^2)^{k/2}|V_n|^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(\beta - \beta_n)^\top V_n^{-1}(\beta - \beta_n) \right\} \\ &\quad \times \frac{b_n^{a_n}}{\Gamma(a_n)} (\sigma^2)^{-a_n-1} \exp \left( -\frac{b_n}{\sigma^2} \right) d\beta \\ &= \frac{b_n^{a_n}}{\Gamma(a_n)} (\sigma^2)^{-a_n-1} \exp \left( -\frac{b_n}{\sigma^2} \right) \\ &\quad \times \int \frac{1}{(2\pi\sigma^2)^{k/2}|V_n|^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(\beta - \beta_n)^\top V_n^{-1}(\beta - \beta_n) \right\} d\beta \\ &= \frac{b_n^{a_n}}{\Gamma(a_n)} (\sigma^2)^{-a_n-1} \exp \left( -\frac{b_n}{\sigma^2} \right) \times 1 \\ &= \frac{b_n^{a_n}}{\Gamma(a_n)} (\sigma^2)^{-a_n-1} \exp \left( -\frac{b_n}{\sigma^2} \right) \end{aligned}$$

La expresión anterior corresponde a la función de densidad de una distribución gamma-inversa con parámetros  $(a_n, b_n)$ . Por lo tanto

$$\sigma^2 | y \sim \text{Inv-Gamma}(a_n, b_n).$$

Por otro lado, la distribución posterior marginal de  $\beta$  se puede encontrar a partir de la

definición, esto es,

$$\begin{aligned}
p(\beta \mid y) &= \int_0^\infty p(\beta, \sigma^2 \mid y) d\sigma^2 \\
&= \int_0^\infty \frac{1}{(2\pi\sigma^2)^{k/2}|V_n|^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\beta - \beta_n)^\top V_n^{-1} (\beta - \beta_n) \right\} \\
&\quad \times \frac{b_n^{a_n}}{\Gamma(a_n)} (\sigma^2)^{-a_n-1} \exp \left( -\frac{b_n}{\sigma^2} \right) d\sigma^2 \\
&= \frac{b_n^{a_n}}{(2\pi)^{k/2}|V_n|^{1/2}\Gamma(a_n)} \int_0^\infty (\sigma^2)^{-(a_n+k/2)-1} \\
&\quad \times \exp \left\{ -\frac{1}{\sigma^2} \left[ b_n + \frac{1}{2} (\beta - \beta_n)^\top V_n^{-1} (\beta - \beta_n) \right] \right\} d\sigma^2 \\
&= \frac{b_n^{a_n}}{(2\pi)^{k/2}|V_n|^{1/2}\Gamma(a_n)} \int_0^\infty (\sigma^2)^{-m-1} \exp \left\{ -\frac{\lambda}{\sigma^2} \right\} d\sigma^2
\end{aligned}$$

donde  $m = a_n + k/2$  y  $\lambda = b_n + \frac{1}{2}(\beta - \beta_n)^\top V_n^{-1}(\beta - \beta_n)$ ,

$$\begin{aligned}
&= \frac{b_n^{a_n}}{(2\pi)^{k/2}|V_n|^{1/2}\Gamma(a_n)} \int_0^\infty (\sigma^2)^{-(m+1)} \exp \left\{ -\frac{\lambda}{\sigma^2} \right\} d\sigma^2 \\
&= \frac{b_n^{a_n}}{(2\pi)^{k/2}|V_n|^{1/2}\Gamma(a_n)} \times \frac{\Gamma(m)}{\lambda^m}
\end{aligned}$$

usando la función de densidad gamma,  $\int_0^\infty x^{-(\alpha+1)} e^{-\beta/x} dx = \frac{\Gamma(\alpha)}{\beta^\alpha}$  con  $\alpha = m$ ,  $\beta = \lambda$ ,

$$\begin{aligned}
&= \frac{b_n^{a_n}}{(2\pi)^{k/2}|V_n|^{1/2}\Gamma(a_n)} \times \frac{\Gamma(a_n + k/2)}{\left[ b_n + \frac{1}{2}(\beta - \beta_n)^\top V_n^{-1}(\beta - \beta_n) \right]^{a_n+k/2}} \\
&= \frac{\Gamma(a_n + k/2)}{\Gamma(a_n)} \frac{b_n^{a_n}}{(2\pi)^{k/2}|V_n|^{1/2}} \left[ b_n + \frac{1}{2}(\beta - \beta_n)^\top V_n^{-1}(\beta - \beta_n) \right]^{-(a_n+k/2)} \\
&= \frac{\Gamma(a_n + k/2)}{\Gamma(a_n)} \frac{1}{(2\pi)^{k/2}|V_n|^{1/2} b_n^{-a_n}} \left[ b_n \left( 1 + \frac{1}{2b_n} (\beta - \beta_n)^\top V_n^{-1} (\beta - \beta_n) \right) \right]^{-(a_n+k/2)} \\
&= \frac{\Gamma(a_n + k/2)}{\Gamma(a_n)} \frac{1}{(2\pi)^{k/2}|V_n|^{1/2}} b_n^{-k/2} \left[ 1 + \frac{1}{2b_n} (\beta - \beta_n)^\top V_n^{-1} (\beta - \beta_n) \right]^{-(a_n+k/2)} \\
&= \frac{\Gamma(\frac{2a_n+k}{2})}{\Gamma(\frac{2a_n}{2}) (2a_n\pi)^{k/2} \left| \frac{b_n}{a_n} V_n \right|^{1/2}} \left[ 1 + \frac{1}{2a_n} (\beta - \beta_n)^\top \left( \frac{b_n}{a_n} V_n \right)^{-1} (\beta - \beta_n) \right]^{-(2a_n+k)/2}
\end{aligned}$$

La expresión anterior corresponde a la función de densidad de una distribución t de Student multivariada con  $2a_n$  grados de libertad, localización  $\beta_n$  y matriz de escala  $\frac{a_n}{b_n} V_n$ .



## Ejercicio

### 3 Conexión con regularización

1. Regresión Ridge: Muestre que si se toma un prior Normal isotrópico

$$\beta \sim \mathcal{N}(0, \tau^2 I)$$

el estimador de máxima a posteriori (MAP) es equivalente a la regresión Ridge:

$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \|y - X\beta\|^2 + \lambda \|\beta\|^2, \quad \lambda = \sigma^2 / \tau^2$$

#### Solución:

Sean  $p(\mathbf{Y} | \boldsymbol{\beta})$  la función de densidad de probabilidad de  $\mathbf{Y}$  dada  $\boldsymbol{\beta}$  la cual es una normal multivariada con media  $X\boldsymbol{\beta}$  y varianza  $\sigma^2 I_n$ , y  $p(\boldsymbol{\beta})$  la función de densidad de probabilidad de  $\boldsymbol{\beta}$  la cual es una normal multivariada con media  $\mathbf{0}$  y varianza  $\tau^2 I_p$ . Observe que si la función de densidad de los datos se le quitan las constantes, se tiene que esta es proporcional a,

$$\begin{aligned} p(\mathbf{Y} | \boldsymbol{\beta}) &\propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} (\mathbf{Y} - X\boldsymbol{\beta})^\top (\mathbf{Y} - X\boldsymbol{\beta})\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} (\mathbf{Y}^\top \mathbf{Y} - 2\boldsymbol{\beta}^\top X^\top \mathbf{Y} + \boldsymbol{\beta}^\top X^\top X\boldsymbol{\beta})\right) \end{aligned}$$

Del mismo modo, se tiene la distribución a priori es proporcional a,

$$p(\boldsymbol{\beta}) \propto \exp\left(-\frac{1}{2\tau^2} \boldsymbol{\beta}^\top \boldsymbol{\beta}\right).$$

De lo anterior se sigue la distribución aposterior es proporcional a,

$$\begin{aligned} p(\boldsymbol{\beta} | \mathbf{Y}) &\propto p(\mathbf{Y} | \boldsymbol{\beta}) \cdot p(\boldsymbol{\beta}) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} (\mathbf{Y}^\top \mathbf{Y} - 2\boldsymbol{\beta}^\top X^\top \mathbf{Y} + \boldsymbol{\beta}^\top X^\top X\boldsymbol{\beta}) - \frac{1}{2\tau^2} \boldsymbol{\beta}^\top \boldsymbol{\beta}\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} \mathbf{Y}^\top \mathbf{Y} + \frac{1}{\sigma^2} \boldsymbol{\beta}^\top X^\top \mathbf{Y} - \frac{1}{2\sigma^2} \boldsymbol{\beta}^\top X^\top X\boldsymbol{\beta} - \frac{1}{2\tau^2} \boldsymbol{\beta}^\top \boldsymbol{\beta}\right) \\ &\propto \exp\left(\boldsymbol{\beta}^\top \left(\frac{X^\top \mathbf{Y}}{\sigma^2}\right) - \frac{1}{2} \boldsymbol{\beta}^\top \left(\frac{X^\top X}{\sigma^2} + \frac{I_p}{\tau^2}\right) \boldsymbol{\beta}\right) \end{aligned}$$

Si se define a

$$\begin{aligned} \Sigma^{-1} &= \frac{X^\top X}{\sigma^2} + \frac{I_p}{\tau^2} \\ \boldsymbol{\mu} &= \Sigma \left( \frac{X^\top \mathbf{Y}}{\sigma^2} \right) \end{aligned}$$

se tiene que la expresión anterior se puede reescribir como,

$$\begin{aligned} p(\boldsymbol{\beta} | \mathbf{Y}) &\propto \exp\left(-\frac{1}{2} \boldsymbol{\beta}^\top \Sigma^{-1} \boldsymbol{\beta} + \boldsymbol{\beta}^\top \Sigma^{-1} \boldsymbol{\mu}\right) \\ &\propto \exp\left(-\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\mu})^\top \Sigma^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu})\right) \end{aligned}$$

...

Observe que la expresión anterior corresponde a el kernel de una distribución normal multivariada con media  $\boldsymbol{\mu}$  y varianza  $\Sigma$ . De aquí, observe que la media explicita de la distribución a posterior de  $\boldsymbol{\beta}$  está dada por,

$$\begin{aligned}\boldsymbol{\mu} &= \left( \frac{X^\top X}{\sigma^2} + \frac{I_p}{\tau^2} \right)^{-1} \frac{X^\top \mathbf{Y}}{\sigma^2} \\ &= \left( X^\top X + \frac{\sigma^2}{\tau^2} I_p \right)^{-1} X^\top \mathbf{Y},\end{aligned}$$

que coincide con el estimador Ridge  $\hat{\boldsymbol{\beta}}$ . Observe que en este caso, la relación entre  $\alpha, \sigma^2$  y  $\tau^2$  es

$$\alpha = \frac{\sigma^2}{\tau^2}$$

la penalización dependerá de la discrepancia que hay entre la varianza a priori y la varianza de los datos.

### Ejercicio

2. Regresión Lasso: Muestre que si en lugar de un prior Normal se utiliza un prior Laplace (doble-exponencial)

$$p(\beta_j) \propto \exp(-\lambda |\beta_j|),$$

el estimador MAP corresponde a la regresión Lasso:

$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \|y - X\beta\|^2 + \lambda \|\beta\|_1$$

### Solución:

En este caso, considere el modelo regresión lineal dado por,  $y = X\beta + \epsilon$  con errores  $\epsilon \sim N(0, \sigma^2 I)$ , la verosimilitud es,

$$p(y|X, \beta) \propto \exp\left(-\frac{1}{2\sigma^2} \|y - X\beta\|^2\right).$$

Además, considere una distribución apriori para  $\beta$  dada por,

$$p(\beta|\lambda) = \prod_{j=1}^p \frac{\lambda}{2} \exp(-\lambda |\beta_j|) \propto \exp(-\lambda \|\beta\|_1).$$

Una vez que se tiene definida la función de verosimilitud y la distribución a priori, se puede obtener la distribución a posterior dada por,

$$\begin{aligned} p(\beta | y) &\propto p(y | \beta) \cdot p(\beta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T \beta)^2}{2\sigma^2}\right) \cdot \prod_{j=1}^p \frac{\lambda}{2} \exp(-\lambda |\beta_j|) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i^T \beta)^2\right) \cdot \exp\left(-\lambda \sum_{j=1}^p |\beta_j|\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \lambda \sum_{j=1}^p |\beta_j|\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \lambda \|\beta\|_1\right) \\ p(\beta|y) &\propto \exp\left(-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \lambda \sum_{j=1}^p |\beta_j|\right). \end{aligned}$$

Por último, observe que el estimador MAP (Maximum A Posteriori) se obtiene maximizando la distribución posterior,

$$\hat{\beta}_{\text{MAP}} = \arg \max_{\beta} p(\beta|y) = \arg \max_{\beta} p(y|\beta)p(\beta),$$

además, note que maximizar  $p(\beta|y)$  es equivalente a minimizar el negativo del log-posterior,

$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \left\{ \frac{1}{2\sigma^2} \|y - X\beta\|^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}.$$

...

Multiplicando por  $2\sigma^2$  y definiendo  $\lambda' = 2\sigma^2\lambda$ ,

$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \left\{ \|y - X\beta\|^2 + \lambda' \|\beta\|_1 \right\}$$

que es exactamente la formulación de la regresión Lasso.



## Ejercicio

### 4 Extensiones: errores no normales

1. Modelos alternativos: Proponga un modelo de regresión donde el error  $\varepsilon$  no siga una distribución Normal. Ejemplos:
  - $\varepsilon \sim \text{Laplace}(0, b)$  (robusto a outliers).
  - $\varepsilon \sim \text{Student- } t(\nu)$  (colas pesadas).
2. Consecuencias metodológicas: Explique cuáles serían las consecuencias sobre:
  - La forma de la verosimilitud.
  - La existencia o no de priors conjugados.
  - Los métodos de inferencia requeridos (MCMC, aproximación variacional, etc.).

### Solución:

Para realizar este análisis, considere un modelo de regresión lineal dado por, dado por,  $y = X\beta + \varepsilon$  con errores  $\varepsilon \sim \text{Student- } t(\nu)$ . En este caso la verosimilitud toma la forma,

$$\begin{aligned} p(y|X, \beta, \sigma^2, \nu) &= \prod_{i=1}^n \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{(y_i - x_i^\top \beta)^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}} \\ &= \left[ \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \right]^n \prod_{i=1}^n \left(1 + \frac{(y_i - x_i^\top \beta)^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}} \end{aligned}$$

Donde,

- $\nu > 0$  son los grados de libertad
- $\sigma^2$  es el parámetro de escala
- $\Gamma(\cdot)$  es la función Gamma
- $x_i^\top$  es la fila  $i$ -ésima de la matriz de diseño  $X$ .

Observe que esta distribución tiene colas más pesadas, lo que lleva a que residuos de mayor tamaño tengan mayor probabilidad, esto produce que los parámetros no se vean tan afectados por outliers y se tengan estimaciones más robustas.

Se sabe que para una verosimilitud de este tipo, no hay distribuciones a priori tales que se pueda tener un modelo conjugado con esta verosimilitud. Sin embargo, se puede recurrir a modelos semi conjugados, los cuales pueden facilitar las cuentas y la implementación de algoritmos de muestreo. Al no poder llegar a una expresión cerrada de la distribución posterior, se debe recurrir a otros métodos de inferencia. Usualmente se utiliza MCMC, como Gibbs sampling o Metropolis Hasting. Particularmente, si se plantea,

$$\varepsilon_i | \lambda_i = l \sim \mathcal{N}(0, \sigma^2/l), \quad \lambda_i \sim \text{Gamma}\left(\frac{\nu}{2}, \frac{\nu}{2}\right)$$

, entonces se puede obtener la distribución t de Student con  $\nu$  grados de libertad y escala  $\sigma^2$ , marginalizando sobre  $\lambda_i$ . A partir de esto, condicionando a  $\boldsymbol{\lambda}, \boldsymbol{\beta}$  sigue una distribución normal.

En este caso en particular, se puede implementar un muestreo de Gibbs, muestreando a  $\beta$  a partir de la normal condicional,  $\sigma^2$  a partir de una gamma inversa (como en el caso usual) y a  $\lambda_i$  a partir de la gamma propuesta,  $\nu$  se puede muestrear a partir de un algoritmo Metrópolis Hasting.