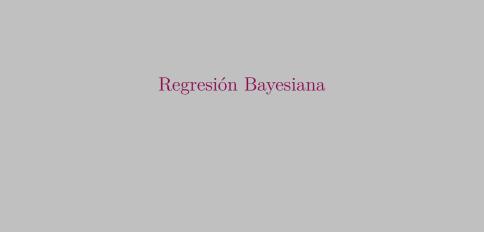
## Estadística Bayesiana Parte II

Irving Gómez Méndez





Regresión Bayesiana

Regresión (Modelo Normal - Inversa Gama)

Detalles sangrientos

Aproximación Normal

Aproximación Normal

#### Tarea 3

Regresión Bayesiana

▶ Dejar tarea 3.



A veces se buscan distribuciones previas que tengan poco impacto en la distribución posterior. Dichas distribuciones son llamadas distribuciones previas de referencia, las cuales son descritas como vagas, constantes, difusas o no informativas.

Aproximación Normal

Por ejemplo, en el caso de la distribución Normal con media  $\theta$  y varianza conocida  $\sigma^2$ , al considerar una previa normal con media  $\mu_0$  y varianza  $\tau_0^2$ , obtuvimos:

$$\theta | \mathbf{Y} \sim \mathsf{Normal}(\mu_n, \tau_n^2),$$

donde

Regresión Bayesiana

$$\mu_n = \frac{\frac{1}{\tau_0^2} \mu_0 + \frac{n}{\sigma^2} \bar{Y}}{\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}} \quad \text{y} \quad \frac{1}{\tau_n^2} = \frac{1}{\tau_0^2} + \frac{n}{\sigma^2}.$$

### Distribuciones impropias

Si hacemos  $\tau_0 \to \infty$ , entonces

$$\theta | \mathbf{Y} \sim \mathsf{Normal}\left( \bar{Y}, \frac{\sigma^2}{n} \right),$$

Aproximación Normal

note que si  $\tau_0 \to \infty$ , entonces  $p(\theta) \propto \mathbb{1}_{\mathbb{R}}(\theta)$ . Lo que va en orden con el principio de razón insuficiente.

Note que la función  $\mathbb{1}_{\mathbb{R}}(\theta)$  no posee integral finita, y por lo tanto no hay manera de normalizarla para que integre 1. Por lo tanto, no hay manera de obtener una densidad y no determina una distribución en sentido estricto. A este tipo de "densidades" que no poseen integral finita se les conoce como densidades impropias.

Retomando el caso binomial, habíamos propuesto  $\theta \sim \text{Uniforme}(0,1)$ . Pero suponga que estamos interesados en  $\phi = -\log \theta$ , si  $\theta \sim \text{Uniforme}(0, 1)$ , entonces  $\phi \sim \text{Exponencial}(1)$ . Lo que viola el principio de razón insuficiente.

Aproximación Normal

Esta ambigüedad en la que no está claro qué debe ser uniforme puede conducir a importantes contradicciones.

Mostrar paradoja de Bertrand.

#### Función score

Regresión Bayesiana

Suponga que  $Y \sim p(Y|\theta_0)$ , definimos la función score como:

$$sc(\theta) = \frac{d}{d\theta} \log p(Y|\theta).$$

$$\mathbb{E}_{Y|\theta_0}[sc(\theta)] = \int_{\mathcal{Y}} \left[ \frac{d}{d\theta} \log p(Y|\theta) \right] p(Y|\theta_0) dY$$
$$= \int_{\mathcal{Y}} \frac{p(Y|\theta_0)}{p(Y|\theta)} \frac{d}{d\theta} p(Y|\theta) dY.$$

Entonces

$$\mathbb{E}_{Y|\theta_0}[sc(\theta_0)] = \int_{\mathcal{V}} \frac{d}{d\theta} p(Y|\theta) \bigg|_{\theta=\theta_0} dY$$

# Condiciones de regularidad

Regresión Bayesiana

Si se pueden intercambiar las operaciones de integración y derivación, entonces

$$\mathbb{E}_{Y|\theta_0}[sc(\theta_0)] = \frac{d}{d\theta} \left[ \int_{\mathcal{Y}} p(Y|\theta) dY \right]_{\theta = \theta_0}$$
$$= \frac{d}{d\theta} 1 \Big|_{\theta = \theta_0}$$
$$= 0.$$

# Información (esperada) de Fisher

Por otro lado, definimos la información esperada de Fisher por unidad muestral como:

$$\mathcal{I}_{\theta_0}(\theta) = \mathbb{E}_{Y|\theta_0}[sc^2(\theta)]$$

$$= \mathbb{E}_{Y|\theta_0} \left[ \frac{1}{p^2(Y|\theta)} \left( \frac{d}{d\theta} p(Y|\theta) \right)^2 \right].$$

Note que

Regresión Bayesiana

$$\frac{d^2}{d\theta^2} \log p(Y|\theta) = \frac{d}{d\theta} \left[ \frac{1}{p(Y|\theta)} \frac{d}{d\theta} p(Y|\theta) \right]$$
$$= -\frac{1}{p^2(Y|\theta)} \left( \frac{d}{d\theta} p(Y|\theta) \right)^2 + \frac{1}{p(Y|\theta)} \frac{d^2}{d\theta^2} p(Y|\theta).$$

Entonces,

Regresión Bayesiana

$$-\mathbb{E}_{Y|\theta_0}\left[\frac{d^2}{d\theta^2}\log p(Y|\theta)\right] = \mathcal{I}_{\theta_0}(\theta) - \mathbb{E}_{Y|\theta_0}\left[\frac{1}{p(Y|\theta)}\frac{d^2}{d\theta^2}p(Y|\theta)\right].$$

Aproximación Normal

Al evaluar en  $\theta_0$  y suponiendo que se pueden intercambiar las operaciones de integración y derivación, obtenemos que

$$\mathcal{I}_{\theta_0}(\theta_0) = -\mathbb{E}_{Y|\theta_0} \left[ \left. \frac{d^2}{d\theta^2} \log p(Y|\theta) \right|_{\theta = \theta_0} \right].$$

Por lo tanto, bajo condiciones de regularidad

$$\mathbb{E}_{Y|\theta_0}[sc(\theta_0)] = 0, \quad \mathbb{V}_{Y|\theta_0}[sc(\theta_0)] = \mathcal{I}_{\theta_0}(\theta_0)$$

Aproximación Normal

У

$$J(\theta_0) \equiv \mathcal{I}_{\theta_0}(\theta_0) = -\mathbb{E}_{Y|\theta_0} \left[ \left. \frac{d^2}{d\theta^2} \log p(Y|\theta) \right|_{\theta=\theta_0} \right].$$

Mostrar ejemplo de caso exponencial.

# **Ejercicios**

1. Sea  $Y|\theta \sim \text{Exponencial}(\theta), p(Y|\theta) = \theta e^{-\theta y} \mathbb{1}_{(0,\infty)}(y).$ Demuestre que la información esperada de Fisher está dada por

Aproximación Normal

$$J(\theta_0) = \frac{1}{\theta_0^2}.$$

2. Sea  $Y|\theta \sim \mathsf{Poisson}(\theta)$ . Demuestre que la información esperada de Fisher está dada por

$$J(\theta_0) = \frac{1}{\theta_0}.$$

3. Sea  $Y|\theta \sim \text{Binomial}(n,\theta)$ . Demuestre que la información esperada de Fisher está dada por

$$J(\theta_0) = \frac{n}{\theta_0(1-\theta_0)}.$$

## Regla de Jeffreys

Regresión Bayesiana

La distribución no informativa de Jeffreys está dada por

$$p(\theta) \propto \sqrt{J(\theta)}$$
.

Aproximación Normal

La idea detrás de esta definición es que cualquier regla para determinar la densidad previa  $p(\theta)$  debería de generar un resultado equivalente al ser aplicada a un parámetro transformado; es decir  $p(\phi)$  calculado a partir de  $p(\theta)$  y el teorema de cambio de variable debería dar el mismo resultado que calculándolo directamente a partir de la información esperada de Fisher para  $\phi$ .

Si  $p(\theta) \propto \sqrt{J(\theta)}$  y  $\phi = \phi(\theta)$  es una transformación 1-1 de  $\theta$ . Entonces  $p(\phi) \propto \sqrt{J(\phi)}$ .

#### Demostración

Usando regla de la cadena, se tiene que

$$\frac{d}{d\phi}\log p(Y|\phi) = \frac{d}{d\theta}\log p(Y|\phi)\frac{d\theta}{d\phi}$$

У

Regresión Bayesiana

$$\frac{d^2}{d\phi^2}\log p(Y|\phi) = \frac{d^2}{d\theta^2}\log p(Y|\phi)\left(\frac{d\theta}{d\phi}\right)^2 + \frac{d}{d\theta}\log p(Y|\phi)\frac{d^2\theta}{d\phi^2}$$

Multiplicando por -1 y tomando el valor esperado:

$$J(\phi) = J(\theta) \left(\frac{d\theta}{d\phi}\right)^2 - \underbrace{\mathbb{E}_{Y|\theta} \left[\frac{d}{d\theta} \log p(Y|\theta)\right]}_{0} \underbrace{\frac{d^2\theta}{d\phi^2}}.$$

Luego

Regresión Bayesiana

$$\sqrt{J(\phi)} = \sqrt{J(\theta)} \left| \frac{d\theta}{d\phi} \right|,$$

es decir

$$p(\phi) \propto p(\theta) \left| \frac{d\theta}{d\phi} \right|.$$

#### Cota de Cramér-Rao

Regresión Bayesiana

Sean  $Y_1, \ldots, Y_n$  variables aleatorias con función de densidad conjunta  $p(\mathbf{Y}|\theta)$ , y sea  $T(\mathbf{Y})$  cualquier función de tal que  $\mathbb{E}_{\mathbf{Y}|\theta}[T(\mathbf{Y})]$  sea una función diferenciable en  $\theta$ . Además, defina la función score (de la muestra) como

Aproximación Normal

$$sc_n(\theta) = \frac{d}{d\theta} \log p(\mathbf{Y}|\theta)$$
$$= \frac{1}{p(\mathbf{Y}|\theta)} \frac{d}{d\theta} p(\mathbf{Y}|\theta),$$

demostramos que, bajo condiciones de regularidad

$$\mathbb{E}_{\mathbf{Y}|\theta}[sc_n(\theta)] = 0.$$

Entonces, bajo dichas condiciones, se satisface que

$$Cov_{\mathbf{Y}|\theta}(T(\mathbf{Y}), sc_n(\theta)) = \mathbb{E}_{\mathbf{Y}|\theta}[T(\mathbf{Y})sc_n(\theta)]$$

$$= \mathbb{E}_{\mathbf{Y}|\theta}\left[T(\mathbf{Y})\frac{1}{p(\mathbf{Y}|\theta)}\frac{d}{d\theta}p(\mathbf{Y}|\theta)\right]$$

$$= \int_{\mathcal{Y}^n} T(\mathbf{y})\frac{d}{d\theta}p(\mathbf{y}|\theta)d\mathbf{y}$$

$$= \frac{d}{d\theta}\int_{\mathcal{Y}^n} T(\mathbf{y})p(\mathbf{y}|\theta)d\mathbf{y}$$

$$= \frac{d}{d\theta}\mathbb{E}_{\mathbf{Y}|\theta}(T(\mathbf{Y})).$$

Aproximación Normal

Regresión Bayesiana

Por otro lado, usando la desigualdad de Cauchy-Schwartz

Aproximación Normal

$$\mathbb{V}_{\mathbf{Y}|\theta}(T(\mathbf{Y}))\mathbb{V}_{\mathbf{Y}|\theta}(sc_n(\theta)) \ge \left(\operatorname{Cov}_{\mathbf{Y}|\theta}(T(\mathbf{Y}), sc_n(\theta))\right)^2$$

$$= \left(\frac{d}{d\theta}\mathbb{E}_{\mathbf{Y}|\theta}(T(\mathbf{Y}))\right)^2$$

$$\Rightarrow \mathbb{V}_{\mathbf{Y}|\theta}(T(\mathbf{Y})) \ge \frac{\left(\frac{d}{d\theta}\mathbb{E}_{\mathbf{Y}|\theta}(T(\mathbf{Y}))\right)^2}{\mathbb{V}_{\mathbf{Y}|\theta}(sc_n(\theta))}.$$

## Relación entre la varianza de un estimador y la información de Fisher

Si,  $Y_1, \ldots, Y_n$  son v.a. i.i.d. con función de densidad  $p(Y|\theta)$ , entonces

Aproximación Normal

$$\mathbb{V}_{\mathbf{Y}|\theta}(T(\mathbf{Y})) \ge \frac{\left(\frac{d}{d\theta} \mathbb{E}_{\mathbf{Y}|\theta}(T(\mathbf{Y}))\right)^2}{nJ(\theta)}.$$

Más aún, si  $T(\mathbf{Y})$  es un estimador insesgado para  $\theta$ , i.e.  $\mathbb{E}_{\mathbf{Y}|\theta}(T(\mathbf{Y})) = \theta$ , entonces

$$\mathbb{V}_{\mathbf{Y}|\theta}(T(\mathbf{Y})) \ge \frac{1}{nJ(\theta)}.$$

Regresión Bayesiana

Aproximación Normal

Por otro lado, esto también indica que no todos los valores de  $\theta$ poseen la misma cantidad de información. De esta manera, entendiendo la distribución previa como una manera de codificar la información que aporta cada valor de  $\theta$  a priori, es que hace sentido tomarla proporcional a la información esperada de Fisher.

Aproximación Normal

Por ejemplo, demostramos que para el caso Binomial  $(Y|\theta \sim \mathsf{Binomial}(n,\theta))$ 

$$J(\theta) = \frac{n}{\theta(1-\theta)},$$

note que  $\theta \to 0$  o  $\theta \to 1$  corresponden a los casos más informativos y es cuando  $J(\theta) \to \infty$ . Mientras que  $J(\theta)$  es mínima cuando  $\theta = 0.5$ .

### Relevancia histórica de la regla de Jeffreys

El problema de la invarianza ante transformaciones monótonas del principio de razón de insuficiencia de Laplace fue una de las mayores críticas a la estadística bayesiana a inicios del siglo XX, realizada, entre otros, por Ronald A. Fisher. Sin embargo, los estudios y aportaciones de Harold Jeffreys sobre previas no informativas que fueran invariantes ante transformaciones volvieron a traer interés sobre el tema.

Aproximación Normal

Regresión Bayesiana

### Relación de la previa de Jeffreys y la divergencia KL

Aproximación Normal

Para el caso uniparamétrico, puede demostrarse que la regla de Jeffreys maximiza la divergencia de Kullback-Leibler. Aunque es posible que Jeffreys ignorara esta afirmación, puede ser que sí tuviera cierta idea de la existencia de una relación entre la regla que propuso y la divergencia de Kullback-Leibler. Sin embargo, cuando existen más parámetros, la regla de Jeffreys no parece ser una buena alternativa (algo que ya había notado el propio Jeffreys).

### Análisis de referencia con pivotales

Si  $U = Y - \theta$  es una variable aleatoria cuya distribución no depende de  $\theta$  ni de Y, entonces U es una cantidad pivotal y  $\theta$  es llamado un parámetro de localización,  $\theta \in \mathbb{R}$ .

Aproximación Normal

Note que

$$\frac{dY}{dU} = 1 \quad \text{y} \quad \left| \frac{d\theta}{dU} \right| = |-1| = 1,$$

luego

$$p(U) = p(Y|\theta) \left| \frac{dY}{dU} \right| = p(Y|\theta)$$

У

$$p(U) = p(\theta|Y) \left| \frac{d\theta}{dU} \right| = p(\theta|Y).$$

Por lo tanto,

$$p(\theta|Y) = p(Y|\theta)$$

y por lo tanto  $p(\theta) \propto \mathbb{1}_{\mathbb{R}}(\theta)$ .

Si  $U = \frac{Y}{A}$  es una variable aletoria cuya distribución no depende de  $\theta$  ni de Y, entonces U es una cantidad pivotal y  $\theta$  es llamado un parámetro de escala,  $\theta > 0$ .

Aproximación Normal

Note que

$$\frac{dY}{dU} = \theta$$

У

$$\frac{d\theta}{dU} = -\frac{Y}{U^2} = -\frac{Y}{Y^2}\theta^2 = -\frac{\theta^2}{Y},$$

luego

$$p(U) = p(Y|\theta) \left| \frac{dY}{dU} \right| = \theta p(Y|\theta)$$

У

$$p(U) = p(\theta|Y) \left| \frac{d\theta}{dU} \right| = \frac{\theta^2}{|y|} p(\theta|Y)$$

$$\frac{\theta^2}{|y|}p(\theta|Y) = \theta p(Y|\theta);$$

Aproximación Normal

entonces

Regresión Bayesiana

$$p(\theta|Y) = \frac{1}{\theta}|y|p(Y|\theta),$$

es decir

$$p(\theta) \propto \frac{1}{\theta} \mathbb{1}_{(0,\infty)}(\theta).$$

### Ejercicio

Demuestre que si  $\theta$  es un parámetro de escala y  $p(\theta) \propto \frac{1}{\theta} \mathbb{1}_{(0,\infty)}(\theta)$ , entonces

$$p(\theta^2) \propto \frac{1}{\theta^2} \mathbb{1}_{(0,\infty)}(\theta^2)$$

У

$$p(\log \theta) \propto \mathbb{1}_{\mathbb{R}}(\log \theta)$$

#### Modelo Uniforme

Sea  $Y|a,b \sim \mathsf{Uniforme}(a,b)$ ,

$$p(Y|a,b) = \frac{1}{b-a} \mathbb{1}_{(a,b)}(Y).$$

Aproximación Normal

Considere la transformación dada por

$$U = \frac{Y - a}{b - a} \Rightarrow Y = (b - a)U + a$$

У

$$\frac{dY}{dU} = b - a,$$

luego

$$p(U|a,b) = \mathbb{1}_{(0,1)}(U).$$

Es decir  $U \sim \mathsf{Uniforme}(0,1)$ , como la distribución de U no depende de a, b ni Y, entonces U es una cantidad pivotal, a es parámetro de localización y b-a es parámetro de escala. Así, podemos proponer la previa no informativa:

Aproximación Normal

$$p(a,b-a) \propto \frac{1}{b-a} \mathbb{1}_{(0,\infty)}(b-a) \mathbb{1}_{\mathbb{R}}(a).$$

Considerando la reparametrización  $\phi(a, b - a) = (a, b)$ , obtenemos la previa no informativa para los parámetros a y b dada por

$$p(a,b) \propto \frac{1}{b-a} \mathbb{1}_{(a,\infty)}(b) \mathbb{1}_{\mathbb{R}}(a).$$

La verosimilitud puede ser escrita como

$$p(Y|a,b) = \frac{1}{b-a} \mathbb{1}_{(-\infty,y)}(a) \mathbb{1}_{(y,\infty)}(b)$$

Aproximación Normal

y la verosimilitud de una muestra observada estaría dada por

$$p(\mathbf{Y}|a,b) = \frac{1}{(b-a)^n} \mathbb{1}_{(-\infty,y_{(1)})}(a) \mathbb{1}_{(y_{(n)},\infty)}(b).$$

Luego, la distribución posterior está dada por

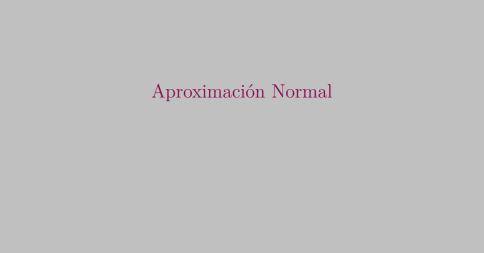
$$p(a, b|\mathbf{Y}) \propto \frac{1}{(b-a)^{n+1}} \mathbb{1}_{(-\infty, y_{(1)})}(a) \mathbb{1}_{(y_{(n)}, \infty)}(b).$$

Al integrar se puede calcular la constante de proporcionalidad y demostrar que la densidad posterior es

Aproximación Normal

$$p(a, b|\mathbf{Y}) = n(n-1) \frac{\left(y_{(n)} - y_{(1)}\right)^{n-1}}{(b-a)^{n+1}} \mathbb{1}_{(-\infty, y_{(1)})}(a) \mathbb{1}_{(y_{(n)}, \infty)}(b)$$

Mostrar modelo Uniforme con previa no informativa.



Suponga que  $Y_1, \ldots, Y_n \stackrel{iid}{\sim} f(Y)$ , pero que nosotros modelamos como  $p(Y|\theta)$ . Además,  $p(\theta)$  es la distribución previa de nuestro modelo. Entonces

$$p(\mathbf{Y}|\theta) = \prod_{i=1}^{n} p(Y_i|\theta)$$

será la verosimilitud de la muestra observada.

Sea  $\theta_0$  el valor que minimiza la divergencia de Kullback-Leibler entre f(Y) y  $p(Y|\theta)$ ,

$$KL(\theta) = \mathbb{E}_{Y \sim f} \left[ \log \left( \frac{f(Y)}{p(Y|\theta)} \right) \right]$$
$$= \int_{\mathcal{Y}} \log \left( \frac{f(Y)}{p(Y|\theta)} \right) f(Y) dY$$

#### Caso discreto

Vamos a demostrar que, cuando n aumenta, la distribución posterior  $p(\theta|\mathbf{Y})$  se concentra alrededor de  $\theta_0$ . Para ello, primero consideraremos el caso en que  $\Theta$  es un espacio discreto.

#### Teorema

Si el espacio parametral  $\Theta$  es finito y  $\mathbb{P}(\theta = \theta_0) > 0$ , entonces  $\mathbb{P}(\theta = \theta_0 | \mathbf{Y}) \xrightarrow[n \to \infty]{} 1$ .

#### Demostración

Considere el logaritmo del cociente de posteriores:

$$\log\left(\frac{p(\theta|\mathbf{Y})}{p(\theta_0|\mathbf{Y})}\right) = \log\left(\frac{p(\theta)}{p(\theta_0)}\right) + \sum_{i=1}^{n} \log\left(\frac{p(Y_i|\theta)}{p(Y_i|\theta_0)}\right)$$

$$\log\left(\frac{p(\theta)}{p(\theta_0)}\right)$$

es una constante, al no depender de n. Por otro lado, note que

$$\mathbb{E}_{Y \sim f} \left[ \log \left( \frac{p(Y|\theta)}{p(Y|\theta_0)} \right) \right]$$

$$= \mathbb{E} \left[ \log f(Y) - \log p(Y|\theta_0) - \log f(Y) + \log p(Y|\theta) \right]$$

$$= KL(\theta_0) - KL(\theta)$$

y por ley fuerte de grandes números se cumple que

$$\frac{1}{n} \sum_{i=1}^{n} \log \left( \frac{p(Y_i | \theta)}{p(Y_i | \theta_0)} \right) \xrightarrow[n \to \infty]{c.s.} KL(\theta_0) - KL(\theta) < 0$$

Por lo tanto,

$$\sum_{i=1}^n \log \left( \frac{p(Y_i|\theta)}{p(Y_i|\theta_0)} \right) \xrightarrow[n \to \infty]{} -\infty,$$

Aproximación Normal

luego

$$\frac{p(\theta|\mathbf{Y})}{p(\theta_0|\mathbf{Y})} \xrightarrow[n \to \infty]{} 0$$

У

$$p(\theta|\mathbf{Y}) \xrightarrow[r \to \infty]{} 0$$
, para todo  $\theta \neq \theta_0$ .

Como la suma de las probabilidades tiene que ser 1, concluimos que

$$p(\theta_0|\mathbf{Y}) \xrightarrow[n\to\infty]{} 1$$

#### Caso continuo

#### Teorema

Sea  $\Theta$  un espacio compacto y A un vecindario de  $\theta_0$  tal que  $\mathbb{P}(\theta \in A) > 0$ , entonces  $\mathbb{P}(\theta \in A|\mathbf{Y}) \longrightarrow 1$ .

Aproximación Normal

#### Demostración

Como  $\Theta$  es compacto, entonces existe una cobertura finita de  $\Theta$ y se puede construir de tal manera que A es el único vecindario que incluye a  $\theta_0$ . Usando el teorema anterior se puede demostrar que la probabilidad posterior para cualquier vecindario que no sea A tiende a 0 cuando  $n \to \infty$  y  $\mathbb{P}(\theta \in A|\mathbf{Y}) \longrightarrow 1$ .

# Convergencia a la distribución Normal usando la información esperada

#### Teorema

Bajo condiciones de regularidad (que incluyen que  $\theta_0$  no esté en la frontera de  $\Theta$ ), la distribución posterior de  $\theta$  es aproximadamente normal con media  $\theta_0$  y varianza  $(nJ(\theta_0))^{-1}$ 

#### Demostración

Sea  $\hat{\theta}$  la moda de la distribución posterior. Luego,

$$\log p(\theta|\mathbf{Y}) = \log p(\hat{\theta}|\mathbf{Y}) + \frac{1}{2}(\theta - \hat{\theta})^2 \left. \frac{d^2}{d\theta^2} \log p(\theta|\mathbf{Y}) \right|_{\theta = \hat{\theta}} + \cdots$$

Note que

$$\left. \frac{d^2}{d\theta^2} \log p(\theta | \mathbf{Y}) \right|_{\theta = \hat{\theta}} = \left. \frac{d^2}{d\theta^2} \log p(\theta) \right|_{\theta = \hat{\theta}} + \sum_{i=1}^n \left. \frac{d^2}{d\theta^2} \log p(Y_i | \theta) \right|_{\theta = \hat{\theta}}$$

Por ley fuerte de grandes números y los teoremas anteriores, tenemos que

$$\frac{1}{n} \sum_{i=1}^{n} \frac{d^2}{d\theta^2} \log p(Y_i|\theta) \bigg|_{\theta = \hat{\theta}} \xrightarrow[n \to \infty]{c.s.} \mathbb{E}_{Y \sim f} \left[ \frac{d^2}{d\theta^2} \log p(Y|\theta) \bigg|_{\theta = \theta_0} \right]$$

Si el modelo de la verosimilitud es correcto, entonces  $f(Y) = p(Y|\theta^*)$  para algún  $\theta^* \in \Theta$ . Y, por lo tanto, la divergencia de Kullback-Leibler se puede escribir como

$$KL(\theta) = \mathbb{E}_{Y \sim f} \left[ \log \left( \frac{p(Y|\theta^*)}{p(Y|\theta)} \right) \right]$$

Recordando que  $KL(\theta) \geq 0$ , podemos verificar fácilmente que  $KL(\theta^*) = 0$  y, por lo tanto,  $\theta_0 = \theta^*$ . Es decir, el parámetro que minimiza la divergencia de Kullback-Leibler es el verdadero parámetro. Entonces

$$\mathbb{E}_{Y \sim f} \left[ \left. \frac{d^2}{d\theta^2} \log p(Y|\theta) \right|_{\theta = \theta_0} \right] = -J(\theta_0)$$

Sabemos que, a medida que  $n \to \infty$ , la distribución se concentra en vecindarios cada vez más pequeños de  $\theta_0$ , y la distancia  $|\hat{\theta} - \theta_0|$  se acerca a cero.

Por lo tanto, al considerar los términos de la serie de Taylor, sólo necesitamos concentrarnos en el término cuadrático, de donde tenemos que

$$p(\theta|\mathbf{Y}) \stackrel{\cdot}{\propto} \exp\left\{-\frac{1}{2}(\theta - \theta_0)^2 (nJ(\theta_0))\right\}$$
$$= \exp\left\{-\frac{(\theta - \theta_0)^2}{2(nJ(\theta_0))^{-1}}\right\}$$

Es decir, para n suficientemente grande,

$$\theta | \mathbf{Y} \stackrel{\cdot}{\sim} \mathsf{Normal}\left(\theta_0, (nJ(\theta_0))^{-1}\right).$$

## Región de $(1-\alpha)$ de probabilidad posterior

Retomando la serie de Taylor, observamos lo siguiente

$$\log p(\theta|\mathbf{Y}) - \log p(\hat{\theta}|\mathbf{Y}) \approx \frac{1}{2}(\theta - \hat{\theta})^2 \left. \frac{d^2}{d\theta^2} \log p(\theta|\mathbf{Y}) \right|_{\theta = \hat{\theta}}$$
$$\Rightarrow -2\log \frac{p(\theta|\mathbf{Y})}{p(\hat{\theta}|\mathbf{Y})} \approx (\theta - \hat{\theta})^2 \left[ -\frac{d^2}{d\theta^2} \log p(\theta|\mathbf{Y}) \right]_{\theta = \hat{\theta}}.$$

Por lo tanto, si  $\theta$  es de dimensión k, entonces:

$$-2\log\frac{p(\theta|\mathbf{Y})}{p(\hat{\theta}|\mathbf{Y})} \stackrel{\cdot}{\sim} \chi_k^2$$

Sea  $q_{\chi_k^2}^{1-\alpha}$  el cuantil de probabilidad  $1-\alpha$  de la distribución  $\chi_k^2$ , i.e.

$$\mathbb{P}\left(\chi_k^2 \le q_{\chi_k^2}^{1-\alpha}\right) = 1 - \alpha.$$

Entonces

$$\mathbb{P}\left[-2\log\frac{p(\theta|\mathbf{Y})}{p(\theta|\mathbf{\hat{Y}})} \le q_{\chi_k^2}^{1-\alpha}\right] \approx 1 - \alpha$$

$$\Rightarrow \mathbb{P}\left[\frac{p(\theta|\mathbf{Y})}{p(\hat{\theta}|\mathbf{Y})} \ge \exp\left\{-\frac{q_{\chi_k^2}^{1-\alpha}}{2}\right\}\right] \approx 1 - \alpha.$$

Es decir, aquella región de  $\Theta$ ,

$$R(\Theta) = \left\{ \theta : \frac{p(\theta|\mathbf{Y})}{p(\hat{\theta}|\mathbf{Y})} \ge \exp\left\{ -\frac{q_{\chi_k^2}^{1-\alpha}}{2} \right\} \right\}$$

Aproximación Normal

corresponde a una región de aproximadamente  $1-\alpha$  de probabilidad posterior.

# Convergencia a la distribución Normal usando la información observada

En el caso de que  $\theta$  sea de dimension k, entonces la serie de Taylor se escribiría como:

$$\log p(\theta|\mathbf{Y}) = \log p(\hat{\theta}|\mathbf{Y}) + \frac{1}{2}(\theta - \hat{\theta})^T \left[ \frac{d^2}{d\theta^2} \log p(\theta|\mathbf{Y}) \right]_{\theta = \hat{\theta}} (\theta - \hat{\theta}) + \cdots$$

у

$$\theta | \mathbf{Y} \stackrel{\cdot}{\sim} \mathsf{Normal}(\hat{\theta}, [I(\hat{\theta})]^{-1}),$$

donde

$$I(\hat{\theta}) = -\left. \frac{d^2}{d\theta^2} \log p(\theta|\mathbf{Y}) \right|_{\theta = \hat{\theta}}$$

► Mostrar ejemplos Beta-Binomial y Gama-Exponencial.

## Aproximación Normal del modelo Beta-Binomial

Suponga que  $Y_1, \ldots, Y_n | \theta \stackrel{iid}{\sim} \mathsf{Bernoulli}(\theta)$ , y considere la previa  $\mathsf{Beta}(\alpha,\beta)$ , entonces  $\theta|\mathbf{Y}\sim\mathsf{Beta}(\alpha^*,\beta^*)$ , donde  $\alpha^* = \alpha + \sum_{i=1}^n y_i \ y \ \beta^* = \beta + n - \sum_{i=1}^n y_i, \text{ luego}$  $p(\theta|\mathbf{Y}) = \text{constante} \times \theta^{\alpha^*-1} (1-\theta)^{\beta^*-1}$ 

У

 $\log p(\theta|\mathbf{Y}) = \text{constante} + (\alpha^* - 1)\log \theta + (\beta^* - 1)\log(1 - \theta).$ 

Ahora calculamos la primera derivada:

$$\frac{d}{d\theta}\log p(\theta|\mathbf{Y}) = \frac{\alpha^* - 1}{\theta} - \frac{\beta^* - 1}{1 - \theta}$$

Aproximación Normal

entonces, la moda posterior de  $\theta$ ,  $\hat{\theta}$ , satisface

$$\frac{1-\hat{\theta}}{\hat{\theta}} = \frac{\beta^* - 1}{\alpha^* - 1}$$

$$\Rightarrow \frac{1}{\hat{\theta}} = \frac{\beta^* - 1}{\alpha^* - 1} + 1$$

$$\Rightarrow \hat{\theta} = \frac{\alpha^* - 1}{\alpha^* + \beta^* - 2}.$$

Calculamos la segunda derivada y la evaluamos en el moda posterior:

$$\left. \frac{d^2}{d\theta^2} \log p(\theta|\mathbf{Y}) \right|_{\hat{\theta}} = -\frac{\alpha^* - 1}{\hat{\theta}^2} - \frac{\beta^* - 1}{(1 - \hat{\theta})^2},$$

note que

$$1 - \hat{\theta} = \hat{\theta} \left( \frac{\beta^* - 1}{\alpha^* - 1} \right),$$

entonces

$$\begin{aligned} \frac{d^2}{d\theta^2} \log p(\theta|\mathbf{Y}) \bigg|_{\hat{\theta}} &= -\frac{\alpha^* - 1}{\hat{\theta}^2} - \frac{(\alpha^* - 1)^2}{\hat{\theta}^2(\beta^* - 1)} \\ &= -\frac{\alpha^* - 1}{\hat{\theta}^2} \left( 1 + \frac{\alpha^* - 1}{\beta^* - 1} \right) \\ &= -\frac{\alpha^* - 1}{\hat{\theta}^2} \left( \frac{\alpha^* + \beta^* - 2}{\beta^* - 1} \right) \\ &= -\frac{\alpha^* + \beta^* - 2}{\hat{\theta}(1 - \hat{\theta})}. \end{aligned}$$

La última igualdad se obtiene notando que

$$\frac{\alpha^* - 1}{\beta^* - 1} = \frac{\hat{\theta}}{1 - \hat{\theta}}.$$

Aproximación Normal

Luego,

Regresión Bayesiana

$$\theta | \mathbf{Y} \stackrel{.}{\sim} \mathcal{N} \left( \hat{\theta}, \frac{\hat{\theta}(1-\hat{\theta})}{\alpha^* + \beta^* - 2} \right),$$

donde

$$\hat{\theta} = \frac{\alpha^* - 1}{\alpha^* + \beta^* - 2}.$$

Note que si hacemos  $\alpha = 1$  y  $\beta = 1$ , entonces  $\hat{\theta} = \bar{y}$  y

$$\theta | \mathbf{Y} \stackrel{\cdot}{\sim} \mathcal{N} \left( \hat{\theta}, \frac{\hat{\theta}(1 - \hat{\theta})}{n} \right)$$

# Aproximación Normal del modelo Normal con previa conjugada

Aproximación Normal

Como se mencionó con anterioridad, si  $Y_1, \ldots, Y_n | \mu, \sigma^2 \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2)$ , entonces la previa conjugada está dada por una distribución Normal-Inversa  $\chi^2$ :

$$\mu|\sigma^2 \sim \mathsf{Normal}(\mu_0, \sigma^2/\kappa_0)$$
 
$$\sigma^2 \sim \mathsf{Inversa} - \chi^2(\nu_0, \sigma_0^2)$$

Se puede demostrar que las distribuciones posteriores están dadas por:

$$\mu | \sigma^2, \mathbf{Y} \sim \mathsf{Normal}(\mu_n, \sigma^2 / \kappa_n)$$
  
 $\sigma^2 | \mathbf{Y} \sim \mathsf{Inversa} - \chi^2(\nu_n, \sigma_n^2),$ 

donde

$$\mu_n = \frac{\kappa_0}{\kappa_0 + n} \mu_0 + \frac{n}{\kappa_0 + n} \bar{y}$$

$$\kappa_n = \kappa_0 + n$$

$$\nu_n = \nu_0 + n$$

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + (n - 1)s^2 + \frac{\kappa_0 n}{\kappa_0 + n} (\bar{y} - \mu_0)^2.$$

Además, la previa de referencia está dada por

$$p(\mu, \sigma^2) \propto \frac{1}{\sigma^2} \mathbb{1}_{\mathbb{R}}(\mu) \mathbb{1}_{(0,\infty)}(\sigma^2).$$

La cual se obtiene fijando los siguientes valores para los hiperparámetros

$$\kappa_0 = 0, \quad \mu_0 \in \mathbb{R}, \quad \nu_0 = -1, \quad \sigma_0^2 = 0,$$

por lo que

$$\kappa_n = n, \quad \mu_n = \bar{y}, \quad \nu_n = n - 1, \quad \sigma_n^2 = s^2,$$

$$\cos s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2.$$

Entonces,

$$p(\mu, \sigma^2 | \mathbf{Y}) = \text{constante} \times (\sigma^2)^{-1/2} \exp \left\{ -\frac{\kappa_n}{2\sigma^2} (\mu - \mu_n)^2 \right\}$$
$$\times (\sigma^2)^{-(\nu_n/2+1)} \exp \left\{ -\frac{\nu_n \sigma_n^2}{2\sigma^2} \right\} \mathbb{1}_{\mathbb{R}}(\mu) \mathbb{1}_{(0,\infty)}(\sigma^2),$$

$$\log p(\mu, \sigma^2 | \mathbf{Y}) = \text{constante} - \left(\frac{\nu_n + 3}{2}\right) \log(\sigma^2) - \frac{\kappa_n}{2\sigma^2} (\mu - \mu_n)^2 - \frac{\nu_n \sigma_n^2}{2\sigma^2},$$

donde reconocemos inmediatamente la moda posterior de  $\mu$ , dada por  $\hat{\mu} = \mu_n$ , y

$$\frac{\partial}{\partial \sigma^2} \log p(\mu, \sigma^2 | \mathbf{Y}) = -\left(\frac{\nu_n + 3}{2}\right) \frac{1}{\sigma^2} + \frac{\kappa_n}{2\sigma^4} (\mu - \mu_n)^2 + \frac{\nu_n \sigma_n^2}{2\sigma^4}.$$

Entonces la moda posterior de  $\sigma^2$ ,  $\hat{\sigma}^2$ , satisface la siguiente ecuación

$$-\left(\frac{\nu_n+3}{2}\right)\frac{1}{\hat{\sigma}^2} + \frac{\nu_n\sigma_n^2}{2\hat{\sigma}^4} = 0,$$

de donde obtenemos

$$\hat{\sigma}^2 = \frac{\nu_n}{\nu_n + 3} \sigma_n^2.$$

#### Ahora calculamos las segundas derivadas

$$\frac{\partial^2}{\partial \mu^2} \log p(\mu, \sigma^2 | \mathbf{Y}) = -\frac{\kappa_n}{\sigma^2},$$

$$\frac{\partial^2}{\partial \mu \partial \sigma} \log p(\mu, \sigma^2 | \mathbf{Y}) = \frac{\kappa_n}{\sigma^4} (\mu - \mu_n),$$

$$\frac{\partial^2}{\partial (\sigma^2)^2} \log p(\mu, \sigma^2 | \mathbf{Y}) = \frac{\nu_n + 3}{2\sigma^4} - \frac{\kappa_n}{\sigma^6} (\mu - \mu_n)^2 - \frac{\nu_n \sigma_n^2}{\sigma^6}.$$

Aproximación Normal

Y evaluamos estas expresiones en las modas posteriores de los parámetros,  $(\hat{\mu}, \hat{\sigma}^2)$ ,

$$\begin{split} \frac{\partial^2}{\partial \mu^2} \log p(\mu, \sigma^2 | \mathbf{Y}) \bigg|_{(\hat{\mu}, \hat{\sigma}^2)} &= -\frac{\kappa_n}{\hat{\sigma}^2}, \\ \frac{\partial^2}{\partial \mu \partial \sigma^2} \log p(\mu, \sigma^2 | \mathbf{Y}) \bigg|_{(\hat{\mu}, \hat{\sigma}^2)} &= 0, \\ \frac{\partial^2}{\partial (\sigma^2)^2} \log p(\mu, \sigma^2 | \mathbf{Y}) \bigg|_{(\hat{\mu}, \hat{\sigma}^2)} &= \frac{\nu_n + 3}{2\hat{\sigma}^4} - \frac{\nu_n \sigma_n^2}{\hat{\sigma}^6} \\ &= \frac{\nu_n + 3}{2\hat{\sigma}^4} - \frac{\nu_n \sigma_n^2}{\hat{\sigma}^4 \nu_n \sigma_n^2} (\nu_n + 3) \\ &= -\frac{\nu_n + 3}{2\hat{\sigma}^4}. \end{split}$$

Entonces,

$$\mu, \sigma^2 | \mathbf{Y} \stackrel{\cdot}{\sim} \mathcal{N} \left( \begin{pmatrix} \hat{\mu} \\ \hat{\sigma}^2 \end{pmatrix}, \begin{pmatrix} \frac{\hat{\sigma}^2}{\kappa_n} & 0 \\ 0 & \frac{2}{\nu_n + 3} \hat{\sigma}^4 \end{pmatrix} \right),$$

donde

$$\hat{\mu} = \mu_n$$
 and  $\hat{\sigma}^2 = \frac{\nu_n}{\nu_n + 3} \sigma_n^2$ .

Si usamos la previa de referencia, obtenemos

$$\mu, \sigma^2 | \mathbf{Y} \stackrel{\cdot}{\sim} \mathcal{N} \left( \begin{pmatrix} \bar{y} \\ \hat{\sigma}^2 \end{pmatrix}, \begin{pmatrix} \frac{\hat{\sigma}^2}{n} & 0 \\ 0 & \frac{2}{n+2} \hat{\sigma}^4 \end{pmatrix} \right),$$

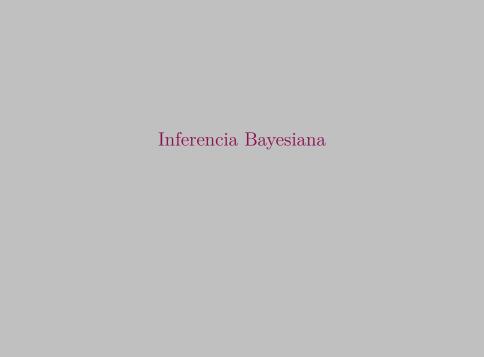
donde

$$\hat{\sigma}^2 = \frac{n-1}{n+2}s^2.$$

Aproximación Normal

#### Tarea 4

▶ Dejar tarea 4.



#### Cómo evitar a Procrustes

Un modelo bayesiano es una máquina que toma de entrada la distribución previa de los parámetros y la verosimilitud y, usando el teorema de Bayes como motor, produce la distribución posterior. Sin embargo, saber la regla matemática del funcionamiento del motor suele ser de muy poca ayuda. Restringirse únicamente a aquellos modelos que permiten la manipulación matemática es una solución procrustea.

Aproximación Normal

Ante este problema es necesario recurrir a alguna técnica numérica que permita aproximar la manipulación matemática.

# Aproximación usando una rendija

Una solución sencilla cuando se tienen pocos parámetros continuos (típicamente uno o dos) consiste en generar una rendija de valores para los parámetros. Sea  $\theta_j$  alguno de estos valores, entonces se puede calcular la distribución posterior en  $\theta_j$  (salvo por una constante de proporcionalidad) usando la fórmula:

$$p(\theta_j|\mathbf{Y}) \propto p(\mathbf{Y}|\theta_j)p(\theta_j).$$

Una importante consecuencia de este hecho es que podemos generar una muestra de la distribución posterior a partir de la rendija de valores propuesta, simplemente basta con seleccionar el valor  $\theta_j$  de manera proporcional a  $p(\mathbf{Y}|\theta_j)p(\theta_j)$ .

## Simular de la predictiva posterior

También nos puede interesar generar una muestra de la distribución predictiva. Una vez que se cuenta con una muestra de la distribución posterior de los parámetros,  $\theta_1, \ldots, \theta_m$ , se puede generar una muestra  $Y_1, \ldots, Y_m$  de la distribución predictiva posterior. Simplemente hay que simular  $Y_i \sim p(Y|\theta_i)$ .

► Mostrar modelo Beta-Binomial.

# Pruebas de hipótesis

Algunas veces la inferencia estadística puede ser formulada como:

- 1. Se cuenta con una hipótesis, la cual puede ser cierta o falsa  $(H: \theta \in \Theta_1)$ .
- 2. Se obtiene evidencia estadística sobre la falsedad de la hipótesis.
- 3. Usamos (o deberíamos usar) el teorema de Bayes para deducir de manera lógica el impacto de la evidencia en la hipótesis

$$\mathbb{P}(H|\mathbf{Y}) = \mathbb{P}(\theta \in \Theta_1|\mathbf{Y}) = \int_{\Theta_1} p(\theta|\mathbf{Y}) d\theta.$$

► Mostrar paradoja de Lindley.

## Estimación por intervalo

Si se cuenta con una muestra de la distribución posterior  $\theta_1, \ldots, \theta_m$ , se puede estimar  $\mathbb{P}(\theta \in \Theta_1 | \mathbf{Y}) = \mathbb{E}\left[\mathbb{1}_{\theta \in \Theta_1} | \mathbf{Y}\right]$  mediante

$$\frac{1}{m} \sum_{j=1}^{m} \mathbb{1}_{\theta_j \in \Theta_1}.$$

También se vuelve sencillo estimar intervalos  $(\theta_1, \theta_2)$  tales que  $\mathbb{P}(\theta \in (\theta_1, \theta_2)) = 1 - \alpha$ . A estos intervalos se les llama intervalos de credibilidad. Un intervalo de particular interés es el de menor longitud cuya probabilidad es  $1 - \alpha$  (highest posterior density interval, HPDI).

## Estimación puntual (MAP)

Recuerde que el estimador bayesiano consiste en toda la distribución posterior. Sin embargo, a veces nos es requerido reportar un único valor. En este caso es común reportar el valor más probable a posteriori (maximum a posteriori, MAP). Lamentablemente, dicho estimador puntual puede dar lugar a resultados absurdos.

### Ejemplo

Considere el ejemplo del globo terráqueo, planteado en la tarea. Suponga que en 3 lanzamientos se obtiene AAA, en este caso el MAP de  $\theta$  vale 1. Lo que es un resultado absurdo.

Aproximación Normal

# Estimación puntual

En vez de reportar el MAP se podría optar por la media o la mediana de la distribución posterior, pero entonces surge la pregunta de qué estimador puntual es el que deberíamos de reportar. Una manera de tomar esta decisión es a través del uso de alguna función de pérdida  $L(\theta, \theta_0)$ , donde  $\theta_0$  es el verdadero parámetro. Pero, jel verdadero parámetro es desconocido!

Aproximación Normal

Para solucionar este inconveniente se minimiza la pérdida esperada, tomando el valor esperado con respecto a la distribución posterior. Es decir se calcula

$$\begin{split} L(\theta) &= \mathbb{E}_{\tilde{\theta} \sim p(\theta|\mathbf{Y})}[L(\theta, \tilde{\theta})] \\ &= \int_{\Theta} L(\theta, \tilde{\theta}) p(\tilde{\theta}|\mathbf{Y}) d\tilde{p} \end{split}$$

y se selecciona  $\hat{\theta} \in \arg\min_{\theta \in \Theta} L(\theta)$ .

Si se cuenta con una muestra  $\theta_1, \dots, \theta_m$  de la distribución posterior, entonces  $L(\theta)$  puede ser estimado mediante

$$\frac{1}{m} \sum_{j=1}^{m} L(\theta, \theta_j)$$

▶ Mostrar modelo Beta-Binomial.

## Cómo determinar qué función de pérdida usar

Considere el caso en que se requiere decidir si ordenar una evacuación o no con base en la velocidad del viento provocado por un huracán.

El riesgo a que haya fallecidos y/o personas afectadas aumenta rápidamente conforme aumenta la rapidez del viento. Sin embargo, también se induce un costo al ordenar una evacuación innecesaria, aunque este es mucho menor.

Por lo tanto, se debería usar una función de pérdida muy asimétrica, que crece rápidamente cuando la velocidad del viento excede nuestra inferencia, pero crece lentamente cuando la velocidad del viento es menor que nuestra inferencia.

## Breve comentario sobre las pruebas de hipótesis

Retomando nuestra discusión sobre las pruebas de hipótesis. De manera más general, lo que se desea es calcular

$$\mathbb{P}(H|\text{evidencia}) = \frac{\mathbb{P}(\text{evidencia}|H)\mathbb{P}(H)}{\mathbb{P}(\text{evidencia})}.$$

Lo más importante es aumentar  $\mathbb{P}(H)$ , lo cual requiere un esfuerzo cognitivo y argumentativo, y no se limita a una simple prueba estadística.