

El método de máxima verosimilitud y la divergencia de Kullback-Leibler

Si f y g son FDP's denotamos por

$$D_{\text{KL}}(f, g) := \int f(x) \log \frac{f(x)}{g(x)} dx$$

a la divergencia de Kullback-Leibler entre f y g .

Si f y g son FDP's denotamos por

$$D_{\text{KL}}(f, g) := \int f(x) \log \frac{f(x)}{g(x)} dx$$

a la divergencia de Kullback-Leibler entre f y g . Recordemos que

- $D_{\text{KL}}(f, g) \geq 0$ para cualesquiera f, g , y que
- $D_{\text{KL}}(f, f) = 0$.

Máxima verosimilitud y máximo a posteriori

Sea $\mathcal{F} = \{f_\theta \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f (con f no necesariamente en \mathcal{F}).

Máxima verosimilitud y máximo a posteriori

Sea $\mathcal{F} = \{f_\theta \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f (con f no necesariamente en \mathcal{F}). Habiendo observado $(X_1, \dots, X_n) = (x_1, \dots, x_n)$, denotamos por

$$L_n(\theta) = L(\theta; x_1, \dots, x_n) = \prod_{i=1}^n f(x_i; \theta)$$

a la función de verosimilitud.

Máxima verosimilitud y máximo a posteriori

Sea $\mathcal{F} = \{f_\theta \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f (con f no necesariamente en \mathcal{F}). Habiendo observado $(X_1, \dots, X_n) = (x_1, \dots, x_n)$, denotamos por

$$L_n(\theta) = L(\theta; x_1, \dots, x_n) = \prod_{i=1}^n f(x_i; \theta)$$

a la función de verosimilitud. Definimos el *estimador de máxima verosimilitud* como

$$\hat{\theta}_n^{\text{EMV}} := \arg \max_{\theta \in \Theta} \{L_n(\theta)\}.$$

Máxima verosimilitud y máximo a posteriori

Sea $\mathcal{F} = \{f_\theta \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f (con f no necesariamente en \mathcal{F}). Habiendo observado $(X_1, \dots, X_n) = (x_1, \dots, x_n)$, denotamos por

$$L_n(\theta) = L(\theta; x_1, \dots, x_n) = \prod_{i=1}^n f(x_i; \theta)$$

a la función de verosimilitud. Definimos el *estimador de máxima verosimilitud* como

$$\hat{\theta}_n^{\text{EMV}} := \arg \max_{\theta \in \Theta} \{L_n(\theta)\}.$$

En el contexto bayesiano suponemos una distribución a priori $p(\theta)$. Definimos el *estimador máximo a posteriori* como

$$\hat{\theta}_n^{\text{MAP}} := \arg \max_{\theta \in \Theta} \{p(\theta | x_1, \dots, x_n)\} = \arg \max_{\theta \in \Theta} \{p(x_1, \dots, x_n | \theta) p(\theta)\}.$$

Máxima verosimilitud y máximo a posteriori

Sea $\mathcal{F} = \{f_\theta \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f (con f no necesariamente en \mathcal{F}). Habiendo observado $(X_1, \dots, X_n) = (x_1, \dots, x_n)$, denotamos por

$$L_n(\theta) = L(\theta; x_1, \dots, x_n) = \prod_{i=1}^n f(x_i; \theta)$$

a la función de verosimilitud. Definimos el *estimador de máxima verosimilitud* como

$$\hat{\theta}_n^{\text{EMV}} := \arg \max_{\theta \in \Theta} \{L_n(\theta)\}.$$

En el contexto bayesiano suponemos una distribución a priori $p(\theta)$. Definimos el *estimador máximo a posteriori* como

$$\hat{\theta}_n^{\text{MAP}} := \arg \max_{\theta \in \Theta} \{p(\theta|x_1, \dots, x_n)\} = \arg \max_{\theta \in \Theta} \{p(x_1, \dots, x_n|\theta)p(\theta)\}.$$

En particular, si θ se distribuye uniformemente (i.e. $p(\theta)$ es constante) y $p(x_1, \dots, x_n|\theta) = \prod_{i=1}^n f(x_i; \theta)$, entonces $\hat{\theta}_n^{\text{MAP}} = \hat{\theta}_n^{\text{EMV}}$.

Por otro lado, denotemos

$$\ell_n(\theta) := \log L_n(\theta) = \log \prod_{i=1}^n f(x_i; \theta) = \sum_{i=1}^n \log f(x_i; \theta).$$

Por otro lado, denotemos

$$\ell_n(\theta) := \log L_n(\theta) = \log \prod_{i=1}^n f(x_i; \theta) = \sum_{i=1}^n \log f(x_i; \theta).$$

Sabemos que maximizar $L_n(\theta)$ es equivalente a maximizar $\ell_n(\theta)$. Más aun, maximizar $\ell_n(\theta)$ es equivalente a maximizar

$$M_n(\theta) = M_n(\theta; x_1, \dots, x_n) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i; \theta)}{f(x_i)}$$

pues $M_n(\theta) = \frac{1}{n}(\ell_n(\theta) - c)$ donde c no depende de θ .

Por otro lado, denotemos

$$\ell_n(\theta) := \log L_n(\theta) = \log \prod_{i=1}^n f(x_i; \theta) = \sum_{i=1}^n \log f(x_i; \theta).$$

Sabemos que maximizar $L_n(\theta)$ es equivalente a maximizar $\ell_n(\theta)$. Más aun, maximizar $\ell_n(\theta)$ es equivalente a maximizar

$$M_n(\theta) = M_n(\theta; x_1, \dots, x_n) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i; \theta)}{f(x_i)}$$

pues $M_n(\theta) = \frac{1}{n}(\ell_n(\theta) - c)$ donde c no depende de θ . Por la LGN,

$$M_n(\theta) \xrightarrow{P} \mathbb{E}_f \left(\log \frac{f(X; \theta)}{f(X)} \right) = \int f(x) \log \frac{f(x; \theta)}{f(x)} dx = -D_{\text{KL}}(f, f_\theta). \quad (1)$$

Por otro lado, denotemos

$$\ell_n(\theta) := \log L_n(\theta) = \log \prod_{i=1}^n f(x_i; \theta) = \sum_{i=1}^n \log f(x_i; \theta).$$

Sabemos que maximizar $L_n(\theta)$ es equivalente a maximizar $\ell_n(\theta)$. Más aun, maximizar $\ell_n(\theta)$ es equivalente a maximizar

$$M_n(\theta) = M_n(\theta; x_1, \dots, x_n) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i; \theta)}{f(x_i)}$$

pues $M_n(\theta) = \frac{1}{n}(\ell_n(\theta) - c)$ donde c no depende de θ . Por la LGN,

$$M_n(\theta) \xrightarrow{P} \mathbb{E}_f \left(\log \frac{f(X; \theta)}{f(X)} \right) = \int f(x) \log \frac{f(x; \theta)}{f(x)} dx = -D_{\text{KL}}(f, f_\theta). \quad (1)$$

Intuitivamente, esperaríamos que de esto se siguiera que

$$\arg \max_{\theta \in \Theta} \{M_n(\theta)\} \xrightarrow{P} \arg \max_{\theta \in \Theta} \{-D_{\text{KL}}(f, f_\theta)\}.$$

Por otro lado, denotemos

$$\ell_n(\theta) := \log L_n(\theta) = \log \prod_{i=1}^n f(x_i; \theta) = \sum_{i=1}^n \log f(x_i; \theta).$$

Sabemos que maximizar $L_n(\theta)$ es equivalente a maximizar $\ell_n(\theta)$. Más aun, maximizar $\ell_n(\theta)$ es equivalente a maximizar

$$M_n(\theta) = M_n(\theta; x_1, \dots, x_n) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i; \theta)}{f(x_i)}$$

pues $M_n(\theta) = \frac{1}{n}(\ell_n(\theta) - c)$ donde c no depende de θ . Por la LGN,

$$M_n(\theta) \xrightarrow{P} \mathbb{E}_f \left(\log \frac{f(X; \theta)}{f(X)} \right) = \int f(x) \log \frac{f(x; \theta)}{f(x)} dx = -D_{\text{KL}}(f, f_\theta). \quad (1)$$

Intuitivamente, esperaríamos que de esto se siguiera que

$$\arg \max_{\theta \in \Theta} \{M_n(\theta)\} \xrightarrow{P} \arg \max_{\theta \in \Theta} \{-D_{\text{KL}}(f, f_\theta)\}.$$

O equivalentemente, que

$$\hat{\theta}_n^{\text{EMV}} \xrightarrow{P} \arg \min_{\theta \in \Theta} \{D_{\text{KL}}(f, f_\theta)\}.$$

Sin embargo, en general esto no es cierto. Por eso pediremos dos cosas:

Sin embargo, en general esto no es cierto. Por eso pediremos dos cosas:

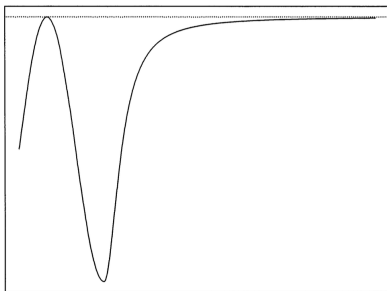
1. Pediremos que $M_n(\theta)$ converga uniformemente en probabilidad a $-D_{\text{KL}}(f, f_\theta)$. (En contraste con $M_n(\theta) \xrightarrow{P} -D_{\text{KL}}(f, f_\theta)$).

Sin embargo, en general esto no es cierto. Por eso pediremos dos cosas:

1. Pediremos que $M_n(\theta)$ converga uniformemente en probabilidad a $-D_{\text{KL}}(f, f_\theta)$. (En contraste con $M_n(\theta) \xrightarrow{P} -D_{\text{KL}}(f, f_\theta)$).
2. Pediremos que la función $\theta \mapsto D_{\text{KL}}(f, f_\theta)$ se maximice en un único punto θ_\star y más aun, pediremos que solo parámetros cercanos θ_\star obtengan valores cercanos a $M(\theta_\star)$.

Sin embargo, en general esto no es cierto. Por eso pediremos dos cosas:

1. Pediremos que $M_n(\theta)$ converga uniformemente en probabilidad a $-D_{\text{KL}}(f, f_\theta)$. (En contraste con $M_n(\theta) \xrightarrow{P} -D_{\text{KL}}(f, f_\theta)$).
2. Pediremos que la función $\theta \mapsto D_{\text{KL}}(f, f_\theta)$ se maximice en un único punto θ_\star y más aun, pediremos que solo parámetros cercanos θ_\star obtengan valores cercanos a $M(\theta_\star)$.



Maximizar $L_n(\theta)$ es asintóticamente equivalente a minimizar D_{KL}

Teorema

Sea $\mathcal{F} = \{f(\cdot; \theta) \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f .

¹En el caso que $f \in \mathcal{F}$, esta hipótesis es consecuencia de otra hipótesis conocida como *identificabilidad* del modelo.

Maximizar $L_n(\theta)$ es asintóticamente equivalente a minimizar D_{KL}

Teorema

Sea $\mathcal{F} = \{f(\cdot; \theta) \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f . Denotemos

$$M_n(\theta) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i; \theta)}{f(x_i)}, \quad M(\theta) = -D_{KL}(f, f_\theta).$$

¹En el caso que $f \in \mathcal{F}$, esta hipótesis es consecuencia de otra hipótesis conocida como *identificabilidad* del modelo.

Maximizar $L_n(\theta)$ es asintóticamente equivalente a minimizar D_{KL}

Teorema

Sea $\mathcal{F} = \{f(\cdot; \theta) \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f . Denotemos

$$M_n(\theta) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i; \theta)}{f(x_i)}, \quad M(\theta) = -D_{KL}(f, f_\theta).$$

Si $\theta_\star := \arg \min_{\theta \in \Theta} \{D_{KL}(f, f_\theta)\}$ esta bien definido¹,

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0, \quad (\text{H1})$$

$$\forall \epsilon > 0, \quad \sup_{\theta: |\theta - \theta_\star| \geq \epsilon} M(\theta) < M(\theta_\star), \quad (\text{H2})$$

¹En el caso que $f \in \mathcal{F}$, esta hipótesis es consecuencia de otra hipótesis conocida como *identificabilidad* del modelo.

Maximizar $L_n(\theta)$ es asintóticamente equivalente a minimizar D_{KL}

Teorema

Sea $\mathcal{F} = \{f(\cdot; \theta) \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f . Denotemos

$$M_n(\theta) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i; \theta)}{f(x_i)}, \quad M(\theta) = -D_{KL}(f, f_\theta).$$

Si $\theta_\star := \arg \min_{\theta \in \Theta} \{D_{KL}(f, f_\theta)\}$ esta bien definido¹,

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0, \quad (\text{H1})$$

$$\forall \epsilon > 0, \quad \sup_{\theta: |\theta - \theta_\star| \geq \epsilon} M(\theta) < M(\theta_\star), \quad (\text{H2})$$

entonces

$$\hat{\theta}_n^{\text{EMV}} \xrightarrow{P} \theta_\star = \arg \min_{\theta \in \Theta} \{D_{KL}(f, f_\theta)\}.$$

¹En el caso que $f \in \mathcal{F}$, esta hipótesis es consecuencia de otra hipótesis conocida como *identificabilidad* del modelo.

Maximizar $L_n(\theta)$ es asintóticamente equivalente a minimizar D_{KL}

Teorema

Sea $\mathcal{F} = \{f(\cdot; \theta) \mid \theta \in \Theta\}$ un modelo paramétrico y sea X_1, \dots, X_n una muestra aleatoria de f . Denotemos

$$M_n(\theta) := \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i; \theta)}{f(x_i)}, \quad M(\theta) = -D_{KL}(f, f_\theta).$$

Si $\theta_\star := \arg \min_{\theta \in \Theta} \{D_{KL}(f, f_\theta)\}$ esta bien definido¹,

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0, \quad (\text{H1})$$

$$\forall \epsilon > 0, \quad \sup_{\theta: |\theta - \theta_\star| \geq \epsilon} M(\theta) < M(\theta_\star), \quad (\text{H2})$$

entonces

$$\hat{\theta}_n^{\text{EMV}} \xrightarrow{P} \theta_\star = \arg \min_{\theta \in \Theta} \{D_{KL}(f, f_\theta)\}.$$

En particular, si $f = f_{\theta_0}$ para algún $\theta_0 \in \Theta$, entonces $\hat{\theta}_n^{\text{EMV}} \xrightarrow{P} \theta_0$.

¹En el caso que $f \in \mathcal{F}$, esta hipótesis es consecuencia de otra hipótesis conocida como *identificabilidad* del modelo.

Demostración. Denotemos $\hat{\theta}_n = \hat{\theta}_n^{\text{EMV}}$.

Demostración. Denotemos $\hat{\theta}_n = \hat{\theta}_n^{\text{EMV}}$. Notemos que

$$\begin{aligned} M(\theta_*) - M(\hat{\theta}_n) &= M(\theta_*) - M_n(\theta_*) + M_n(\theta_*) - M(\hat{\theta}_n) \\ &\leq \sup_{\theta \in \Theta} |M(\theta) - M_n(\theta)| + M_n(\hat{\theta}_n) - M(\hat{\theta}_n) \\ &\quad \quad \quad (\hat{\theta}_n \text{ maximiza } M_n) \\ &\xrightarrow{P} 0. \quad \quad \quad (\text{por (1) y (H1)}) \end{aligned}$$

Demostración. Denotemos $\hat{\theta}_n = \hat{\theta}_n^{\text{EMV}}$. Notemos que

$$\begin{aligned} M(\theta_*) - M(\hat{\theta}_n) &= M(\theta_*) - M_n(\theta_*) + M_n(\theta_*) - M(\hat{\theta}_n) \\ &\leq \sup_{\theta \in \Theta} |M(\theta) - M_n(\theta)| + M_n(\hat{\theta}_n) - M(\hat{\theta}_n) \\ &\xrightarrow{P} 0. \end{aligned} \quad \begin{array}{l} (\hat{\theta}_n \text{ maximiza } M_n) \\ (\text{por (1) y (H1)}) \end{array}$$

En particular, para cualquier $\delta > 0$

$$P \left(M(\hat{\theta}_n) < M(\theta_*) - \delta \right) \rightarrow 0. \quad (2)$$

Demostración. Denotemos $\hat{\theta}_n = \hat{\theta}_n^{\text{EMV}}$. Notemos que

$$\begin{aligned} M(\theta_*) - M(\hat{\theta}_n) &= M(\theta_*) - M_n(\theta_*) + M_n(\theta_*) - M(\hat{\theta}_n) \\ &\leq \sup_{\theta \in \Theta} |M(\theta) - M_n(\theta)| + M_n(\hat{\theta}_n) - M(\hat{\theta}_n) \\ &\quad (\hat{\theta}_n \text{ maximiza } M_n) \\ &\xrightarrow{P} 0. \quad (\text{por (1) y (H1)}) \end{aligned}$$

En particular, para cualquier $\delta > 0$

$$P \left(M(\hat{\theta}_n) < M(\theta_*) - \delta \right) \rightarrow 0. \quad (2)$$

Por otro lado, sea $\epsilon > 0$ arbitrario y sea

$$\delta = M(\theta_*) - \sup_{\theta: |\theta - \theta_*| \geq \epsilon} M(\theta) > 0. \quad (\text{por (H2)})$$

Demostración. Denotemos $\hat{\theta}_n = \hat{\theta}_n^{\text{EMV}}$. Notemos que

$$\begin{aligned} M(\theta_*) - M(\hat{\theta}_n) &= M(\theta_*) - M_n(\theta_*) + M_n(\theta_*) - M(\hat{\theta}_n) \\ &\leq \sup_{\theta \in \Theta} |M(\theta) - M_n(\theta)| + M_n(\hat{\theta}_n) - M(\hat{\theta}_n) \\ &\quad (\hat{\theta}_n \text{ maximiza } M_n) \\ &\xrightarrow{P} 0. \quad (\text{por (1) y (H1)}) \end{aligned}$$

En particular, para cualquier $\delta > 0$

$$P \left(M(\hat{\theta}_n) < M(\theta_*) - \delta \right) \rightarrow 0. \quad (2)$$

Por otro lado, sea $\epsilon > 0$ arbitrario y sea

$$\delta = M(\theta_*) - \sup_{\theta: |\theta - \theta_*| \geq \epsilon} M(\theta) > 0. \quad (\text{por (H2)})$$

De esta manera, $|\theta - \theta_*| \geq \epsilon$ implica $M(\theta) < M(\theta_*) - \delta$.

Demostración. Denotemos $\hat{\theta}_n = \hat{\theta}_n^{\text{EMV}}$. Notemos que

$$\begin{aligned} M(\theta_*) - M(\hat{\theta}_n) &= M(\theta_*) - M_n(\theta_*) + M_n(\theta_*) - M(\hat{\theta}_n) \\ &\leq \sup_{\theta \in \Theta} |M(\theta) - M_n(\theta)| + M_n(\hat{\theta}_n) - M(\hat{\theta}_n) \\ &\quad (\hat{\theta}_n \text{ maximiza } M_n) \\ &\xrightarrow{P} 0. \quad (\text{por (1) y (H1)}) \end{aligned}$$

En particular, para cualquier $\delta > 0$

$$P \left(M(\hat{\theta}_n) < M(\theta_*) - \delta \right) \rightarrow 0. \quad (2)$$

Por otro lado, sea $\epsilon > 0$ arbitrario y sea

$$\delta = M(\theta_*) - \sup_{\theta: |\theta - \theta_*| \geq \epsilon} M(\theta) > 0. \quad (\text{por (H2)})$$

De esta manera, $|\theta - \theta_*| \geq \epsilon$ implica $M(\theta) < M(\theta_*) - \delta$. Por lo tanto,

$$P \left(|\hat{\theta}_n - \theta_*| > \epsilon \right) \leq P \left(M(\hat{\theta}_n) < M(\theta_*) - \delta \right) \xrightarrow{P} 0.$$

□