

ÍNDICE

1.	INTRODUCCIÓN A LA PROBABILIDAD	11
1.1.	Probabilidad, espacio muestral y sucesos	11
1.1.1.	Espacio muestral y sucesos	11
1.1.2.	Probabilidad	14
1.1.3.	Variable aleatoria y función de distribución	15
1.2.	Valores esperados y momentos de una distribución	17
1.3.	Distribuciones conjuntas	19
1.3.1.	Distribuciones marginales	19
1.3.2.	Momentos de una distribución conjunta ..	20
1.3.3.	Distribuciones condicionadas	21
1.4.	Muestra aleatoria simple	23
1.5.	Ejercicios	23
2.	DISTRIBUCIONES	27
2.1.	Distribuciones discretas	27
2.1.1.	Uniforme	27
2.1.2.	Bernoulli	29
2.1.3.	Binomial	30
2.1.4.	Multinomial	32

2.1.5. Poisson	34
2.2. Distribuciones continuas	37
2.2.1. Uniforme	37
2.2.2. Normal.....	38
2.2.3. Normal multivariante	40
2.2.4. Exponencial.....	42
2.3. Ejercicios	44
3. ESTIMACIÓN	47
3.1. Estimación por el método de los momentos	47
3.2. Estimación por mínimos cuadrados.....	49
3.3. Estimación por máxima verosimilitud	50
3.4. Estimación bayesiana.....	53
3.5. Propiedades de los estimadores puntuales	55
3.6. Propiedades asintóticas de los estimadores má- ximo verosímiles	57
3.7. Ejercicios	59
4. CONTRASTE DE HIPÓTESIS.....	61
4.1. Conceptos fundamentales	61
4.2. Contraste mediante intervalos de confianza	64
4.3. Teorema de Neyman-Pearson	66
4.4. Test de la razón de verosimilitud generalizado...	68
4.5. Test q de Pearson	71
4.6. Ejercicios	74
5. MODELOS LINEALES.....	77
5.1. Regresión lineal	77
5.1.1. Modelo estadístico	77
5.1.2. Estimación	79
5.2. Análisis factorial.....	85
5.2.1. Modelo estadístico	86
5.2.2. Estimación	89

ÍNDICE	9
5.2.3. Rotaciones	92
5.2.4. Casos heywood.....	93
5.2.5. Test del modelo	93
5.3. Ejercicios	94
6. TABLAS DE CONTINGENCIA	97
6.1. Esquemas de muestreo.....	98
6.2. Contraste de hipótesis	101
6.3. Introducción a los modelos log-lineales	103
6.4. Ejercicios	106
7. ESTADÍSTICA BAYESIANA	109
7.1. Distribuciones previa y posterior.....	109
7.2. Estimación de parámetros	114
7.3. Evaluación del modelo.....	116
7.4. Comparación de modelos	120
7.5. Algunas consideraciones sobre los métodos esta- dísticos.....	124
7.5.1. Estimación	124
7.5.2. Contraste	125
7.6. Ejercicios	126
BIBLIOGRAFÍA COMENTADA	129
ÍNDICE DE MATERIAS	133

CAPÍTULO 3 ESTIMACIÓN

El problema de la estimación consiste en realizar una inferencia sobre el valor de los parámetros de un modelo estadístico a partir de los datos contenidos en una muestra. Existen tres métodos fundamentales: método de los momentos, mínimos cuadrados y máxima verosimilitud. Este curso trata fundamentalmente del método de máxima verosimilitud.

En general, un parámetro cualquiera se va a designar por θ y su correspondiente estimador por $\hat{\theta}$. El estimador $\hat{\theta}$ es una función de los datos muestrales, cuyo valor depende de la muestra concreta. Por tanto, el parámetro θ toma un valor fijo mientras que el estimador $\hat{\theta}$ es aleatorio, su valor depende de la muestra.

3.1. Estimación por el método de los momentos

Supongamos que se desean estimar k parámetros $(\theta_1, \theta_2, \dots, \theta_k)$ de una distribución $f(X; \theta_1, \theta_2, \dots, \theta_k)$. El método de los momentos consiste en buscar los valores de los parámetros que igualan los k primeros momentos con relación al origen en la muestra y en la población. Por tanto, no se utiliza la distribución completa para obtener estimadores sino únicamente los momentos indicados.

En la población el momento de orden i se ha definido:

$$\alpha_i = \int_{-\infty}^{\infty} X^i f(X) dX$$

el correspondiente momento muestral es:

$$a_i = \frac{\sum_{i=1}^n X^i}{n}$$

De esta forma se obtiene el siguiente sistema de k ecuaciones con k incógnitas:

$$\left. \begin{array}{l} \alpha_1 = a_1 \\ \alpha_2 = a_2 \\ \vdots \\ \alpha_k = a_k \end{array} \right\}$$

Los valores de los parámetros que satisfacen el sistema de ecuaciones constituyen los estimadores $(\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)$.

Ejemplo 3.1. (Distribución de Poisson).

Se obtiene una muestra X_1, \dots, X_n de una distribución de Poisson. Estimar λ por el método de los momentos.

En esta distribución sólo hay un parámetro que además cumple $\alpha = E(X) = \lambda$. El primer momento muestral es $a_1 = \bar{X}$. El estimador es el valor que resuelve: $\alpha_1 = a_1$, es decir $\hat{\lambda} = \bar{X}$.

Ejemplo 3.2. (Regresión logística).

Sea Y una variable dicotómica cuya distribución depende de una variable independiente X . Se desea estimar el parámetro β de la regresión logística:

$$f(Y_i | X) = \frac{\exp(Y\beta X)}{1 + \exp(\beta X)}$$

Dado que Y es dicotómica, para cada valor de X el primer momento poblacional es su valor esperado:

$$E(Y|X) = 0 \times f(Y=0|X) + 1 \times f(Y=1|X) = f(Y=1|X)$$

El correspondiente momento muestral es $\bar{Y}(X)$, que representa la media de Y para todas las observaciones en las que X toma un valor determinado. Por tanto, si X toma los valores x_1, x_2, \dots, x_n para cada uno de ellos se dispone de un momento poblacional: $E(Y|x_1), E(Y|x_2), \dots, E(Y|x_n)$ y un momento muestral: $\bar{Y}(x_1), \bar{Y}(x_2), \dots, \bar{Y}(x_n)$. El estimador de β es el valor que iguala los momentos en la población y la muestra, es decir:

$$\sum_{i=1}^n \frac{\exp(\beta x_i)}{1 + \exp(\beta x_i)} = \sum_{i=1}^n \bar{Y}(x_i)$$

En esta ecuación no es posible despejar β . Para resolverla pueden utilizarse los denominados métodos numéricos, por ejemplo el de la bisección, Newton-Raphson, etc.

3.2. Estimación por mínimos cuadrados

Consiste en asignar a los parámetros aquel valor que minimice la diferencia al cuadrado entre los datos observados y los predichos por el modelo estadístico.

El método puede entenderse con un ejemplo. Supongamos que Y es una variable aleatoria. Se desea estimar la regresión lineal de Y sobre un predictor X , es decir para un sujeto i : $Y_i = \alpha + \beta X_i + E_i$, donde E_i es un error aleatorio con media 0 y varianza σ^2 para todos los valores de X .

Los parámetros del modelo son α y β . Para cada valor de X el valor predicho de Y es: $E(Y) = \alpha + \beta X$, por tanto para estimar los parámetros se minimiza la diferencia cuadrática entre los datos observados y las predicciones del modelo. Si el tamaño muestral es n la diferencia es:

$$d(\alpha, \beta) = \sum_{i=1}^n (Y_i - \alpha - \beta X_i)^2$$

Los estimadores son aquellos valores que minimizan la función $d(\alpha, \beta)$, es decir aquellos valores que hagan que su primera derivada sea 0:

$$\frac{\partial}{\partial \alpha} d(\alpha, \beta) = -2 \sum_{i=1}^n (Y_i - \alpha - \beta X_i) = 0$$

$$\frac{\partial}{\partial \beta} f(\alpha, \beta) = -2 \sum_{i=1}^n (Y_i - \alpha - \beta X_i) X_i = 0$$

Desarrollando estas ecuaciones se encuentran los estimadores descritos en cursos anteriores (Pardo y San Martín, 1999): $\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$ y $\hat{\beta} = \text{Cov}(X, Y) / \text{Std}(X) \text{Std}(Y)$.

3.3. Estimación por máxima verosimilitud

El método de máxima verosimilitud consiste en asignar a los parámetros aquel valor que haga máxima la probabilidad de los datos observados. A diferencia de los anteriores utiliza la distribución completa de probabilidad.

En caso de que la distribución de la variable sea normal los estimadores de mínimos cuadrados y máxima verosimilitud coinciden. Continuando con el ejemplo anterior, la distribución de Y es:

$$f(Y_i) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{Y_i - \alpha - \beta X_i}{\sigma}\right)^2\right)$$

Por tanto la función de densidad de la muestra compuesta por n observaciones es:

$$L(\alpha, \beta) = \prod_{i=1}^n f(Y_i) = \frac{1}{\sigma^n (2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{(Y_i - \alpha - \beta X_i)^2}{\sigma^2}\right)$$

Los estimadores máxima verosímiles son los valores que maximizan la función de verosimilitud $L(\alpha, \beta)$, es decir la función de densidad de la muestra. Analizando la forma de esta función

puede verse que estos valores son aquellos que minimizan el término en el exponente:

$$\sum_{i=1}^n \frac{(Y_i - \alpha - \beta X_i)^2}{\sigma^2}$$

La forma de este término es la misma que la de la función $d(\alpha, \beta)$ comentada en relación con la estimación por mínimos cuadrados. Por tanto, los estimadores obtenidos por ambos métodos coinciden. Esto es una consecuencia de haber asumido que la distribución de Y es normal. Con cualquier otra distribución la forma de proceder sería la misma, obtener la función de verosimilitud y maximizarla con respecto a los parámetros.

Ejemplo 3.3. (Distribución de Bernoulli)

Supongamos que a un sujeto se le presenta 25 veces una determinada tarea. El resultado de cada presentación se clasifica como éxito o fracaso y se considera que la probabilidad de éxito \neq no cambia a lo largo de las presentaciones. ¿Cual es la probabilidad de éxito asumiendo independencia entre distintas presentaciones?

La variable X_i describe el resultado de la presentación i y sigue la distribución de Bernoulli:

$$f(X_i; \pi) = \pi^{X_i} (1 \pm \pi)^{(1 \pm X_i)}$$

La función de probabilidad del vector de resultados de las 25 presentaciones tiene la forma:

$$f(X; \pi) = \prod_{i=1}^{25} \pi^{X_i} (1 \pm \pi)^{(1 \pm X_i)}$$

Supongamos que el numero de éxitos se indica por z , siendo $z = \sum_{i=1}^{25} X_i$. Entonces la función de verosimilitud es:

$$L(\pi) = \pi^z (1 \pm \pi)^{(25 \pm z)}$$

En la práctica se trabaja con el logaritmo de la función de verosimilitud por mayor sencillez matemática:

$$\log L(\pi) = z \log \pi + (25 - z) \log(1 - \pi)$$

El estimador máximo verosímil es el valor que maximiza $\log L(\pi)$. Un criterio necesario, pero no suficiente, para que a un valor concreto de π le corresponda un máximo de $\log L(\pi)$ es que anule su primera derivada. En el ejemplo:

$$\begin{aligned} \frac{\partial}{\partial \pi} \log L(\pi) &= \frac{z}{\pi} - \frac{25 - z}{1 - \pi} \\ &= 0 \end{aligned}$$

La ecuación anterior se denomina ecuación de estimación. Su solución es el estimador máximo verosímil de π :

$$\begin{aligned} \frac{z}{n} &= \frac{25 - z}{1 - \pi} \\ z - z\pi &= 25\pi - z\pi \\ \hat{\pi} &= \frac{z}{25} \end{aligned}$$

En realidad el criterio de la primera derivada no es suficiente para determinar que $\hat{\pi}$ es un estimador máximo verosímil porque dicha derivada se anula tanto si la función tiene un máximo en $\hat{\pi}$ cómo si tiene un mínimo. En caso de que la función tenga un máximo se cumple que su segunda derivada es negativa. En el ejemplo:

$$\frac{\partial}{\partial \pi^2} \log L(\pi) = -\frac{z}{\pi^2} - \frac{25 - z}{(1 - \pi)^2}$$

La cual es necesariamente menor que 0, por lo que efectivamente $\hat{\pi}$ es un máximo de $\log L(\pi)$.

Ejemplo 3.4.

Continuando con el mismo ejemplo, supongamos que se relaja el supuesto de independencia y se asume únicamente independencia condicional. Supongamos que la probabilidad de éxito en la tarea i depende del número de aciertos y errores cometidos en las $i - 1$ tareas anteriores. El número de aciertos anteriores es $C_i = \sum_{j=1}^{i-1} X_j$ y el de errores $D_i = \sum_{j=1}^{i-1} (1 - X_j)$. Podría plantearse el modelo:

$$\pi_i = \frac{\exp(\alpha + \beta C_i + \delta D_i)}{1 + \exp(\alpha + \beta C_i + \delta D_i)}$$

La cantidad π_i debe interpretarse como la probabilidad del ensayo i condicionada en los ensayos anteriores. La función de verosimilitud se obtiene a partir de la regla del producto y tiene la forma:

$$L(\alpha, \beta, \delta) = \prod_{i=1}^{25} \pi_i^{X_i} (1 - \pi_i)^{(1 - X_i)}$$

3.4. Estimación bayesiana

La estimación bayesiana se caracteriza porque permite incorporar las expectativas del investigador acerca de los valores de los parámetros. La estadística bayesiana se diferencia la denominada estadística frecuentista en dos aspectos importantes:

- En la estadística bayesiana tanto los datos como los parámetros son cantidades aleatorias. En los métodos anteriores los parámetros se consideran cantidades fijas.
- Al ser los parámetros cantidades aleatorias siguen una distribución, denominada *distribución previa*, que expresa las expectativas del investigador en ausencia de ningún dato.

En los métodos frecuentistas el valor que toman los estimadores depende únicamente de cuales hayan sido los datos ob-

servados. Por el contrario, en los métodos bayesianos dichas valores dependen de los datos y también de las expectativas previas. Por esta razón existen distintas escuelas en el seno de la estadística, dependiendo de sí se considera legítimo o no utilizar distribuciones previas. Lo cierto es que desde un punto de vista aplicado los métodos bayesianos pueden ser de gran utilidad para determinados problemas. Además, a medida que el tamaño muestral aumenta los estimadores dependen más de los datos y menos de las distribuciones previas, por lo que en el límite $n \rightarrow \infty$ el estimador máximo verosímil y el bayesiano coinciden.

A modo de ejemplo, consideremos el problema de estimar la media de una distribución normal $f(Y|\mu)$ con varianza conocida a partir de n réplicas de un experimento aleatorio. El estimador máximo verosímil es aquel que maximiza la función de verosimilitud:

$$f(Y|\mu) = \prod_{i=1}^n f(Y_i|\mu)$$

Supongamos que la variable en cuestión mide el cociente intelectual de los sujetos. Se sabe que en la población general su distribución es normal (100, 16). Por tanto, el investigador puede suponer que en la subpoblación con la que está trabajando el valor de la media oscilará en torno a 100, siendo muy inesperado encontrar medias, digamos, superiores a 150 o inferiores a 50. Esta expectativa se formaliza mediante la definición de una distribución previa. Por ejemplo, la distribución previa de μ puede ser la normal (100, 20), denominada $f(\mu)$. Es importante advertir que esta distribución previa no depende de ningún dato observado en el experimento actual sino que se define arbitrariamente.

En la estadística bayesiana un objetivo es modificar las expectativas previas del investigador de acuerdo con la evidencia encontrada en la muestra. Según el teorema de Bayes la distribución posterior de μ puede obtenerse del siguiente

modo, a partir de la función de verosimilitud y la distribución previa:

$$f(\mu | Y) = \frac{f(Y | \mu)f(\mu)}{f(Y)}$$

Siendo $f(\mathbf{Y})$ la distribución marginal de \mathbf{Y} , es decir: $f(\mathbf{Y}) = \int f(\mathbf{Y} | \mu) f(\mu) d\mu$. El estimador bayesiano de μ es la distribución posterior $f(\mu | \mathbf{Y})$. Una descripción mas completa de la estadística bayesiana aparece en el capítulo 7.

3.5. Propiedades de los estimadores puntuales

Hasta el momento se han descrito varios métodos mas o menos intuitivos para obtener estimadores. Puede plantearse la cuestión de en qué sentido son unos mejores que otros. Para verificarlo se han descrito una serie de propiedades que deberían cumplir los estimadores. Para un valor fijo de θ se define la distribución muestral del estimador $\hat{\theta}$, de la cual dependen las propiedades de dicho estimador:

1. *Insesgado*. Un estimador es insesgado cuando cumple:

$$E(\hat{\theta}) = \theta$$

para todo valor de θ en el espacio paramétrico. A la cantidad $b(\hat{\theta}) = E(\hat{\theta}) - \theta$ se le denomina sesgo del estimador.

Ejemplo 3.5.

Según se ha estudiado en cursos anteriores, para estimar la media poblacional μ se utiliza la media muestral: $\bar{X} = \sum_{i=1}^n X_i/n$. Se trata de un estimador insesgado dado que $E(\bar{X}) = \mu$.

2. *Consistente*. Un estimador es consistente si para cualquier $\epsilon > 0$:

$$\lim_{n \rightarrow \infty} P(|\theta - \hat{\theta}| > \epsilon) = 0$$

Un estimador insesgado y que cumpla $Var(\hat{\theta}) \rightarrow 0$ cuando $n \rightarrow \infty$ se dice que es consistente.

Ejemplo 3.6.

La varianza de la media muestral es σ^2/n . Esta cantidad tiende a 0 por lo que la media muestral es un estimador consistente de μ .

3. *Eficiente*. Un estimador es eficiente en caso de que sea insesgado y su varianza mínima. Para estimar un parámetro pueden utilizarse distintos estimadores con diferente varianza. Por tanto es deseable utilizar aquel que tenga menos dispersión.

Las diferencias entre el parámetro y el estimador obtenido en las posibles muestras pueden cuantificarse mediante el error cuadrático medio, cuyo valor depende de las tres propiedades anteriores. Se conoce como $MSE(\hat{\theta})$ (*Mean squared error o error cuadrático medio*) a la cantidad $MSE(\hat{\theta}) = E((\hat{\theta} - \theta)^2)$. Desarrollando esta expresión se llega a:

$$\begin{aligned} MSE(\hat{\theta}) &= E((\hat{\theta} - \theta)^2) \\ &= E(([\hat{\theta} - E(\hat{\theta})] - [E(\hat{\theta}) - \theta])^2) \\ &= E((\hat{\theta} - E(\hat{\theta}))^2) + E((E(\hat{\theta}) - \theta)^2) - 2E((\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta)) \\ &= Var(\hat{\theta}) + b(\hat{\theta})^2 \end{aligned}$$

El MSE es la suma de la varianza del estimador más su sesgo al cuadrado y muestra la relación entre estas tres cantidades.