

Fundamentos



Teoría Estadística

Índice general

1. Estimación Puntual	1
1.1. Métodos	1
1.1.1. Método Pivotal	1
2. Intervalos para medias y varianzas poblacionales	3
2.1. Métodos	3
2.1.1. Métodos para la Media	3
2.1.2. Métodos para la Varianza	3
2.1.3. Métodos para comparación de medias	4
2.1.4. Método para la Comparación de Varianzas	4
2.1.5. Métodos para una proporción y diferencia de proporciones	5
2.2. Practica	5
3. Contraste Hipótesis estadísticas	7
3.1. Preliminares	7
3.1.1. Hipótesis estad	7
3.1.2. Tipos de Error	7
3.1.3. Estad	8
3.2. Procedimientos de contraste de hipótesis	9
3.2.1. Pruebas comunes con muestras grandes	9
3.2.2. Pruebas para muestras pequeñas	11
3.2.3. Niveles de significancia y p-values	11
4. Lemma Neyman Pearson	12
4.1. Concepto de Espacio paramétrico y potencia	12
4.2. Neyman-Pearson	12
5. Estadística Bayesiana	13
5.1. Contraste con el Enfoque frecuentista	13

Lista de Ecuaciones

1. Estimación Puntual	1
1.1. Intervalo de Confianza Bilateral	1
1.2. Intervalo de Confianza Unilateral	1
2. Intervalos para medias y varianzas poblacionales	3
2.2. Intervalo de Confianza Media con σ desconocida	3
2.4. Intervalo de Confianza Varianza	3

2.6.	Intervalo de Confianza Comparación de Medias.	4
2.8.	Intervalo de Confianza para comparación de Varianzas	4
2.9.	Intervalo de Confianza para una Proporción	5
2.10.	Intervalo de Confianza para comparación de Proporciones	5
3.	Contraste Hipótesis estadísticas	7
3.4.	Pruebas de Hipótesis de nivel I para muestras grandes	9
3.8.	Pruebas de Hipótesis de nivel I para proporciones	10
3.11.	Pruebas de Hipótesis de nivel I para comparar proporciones	10
4.	Lemma Neyman Pearson	12
4.1.	Funcion de potencia	12
5.	Estadística Bayesiana	13
5.1.	Teorema de Bayes	13
5.2.	Teorema de Bayes distribución de Theta	13

Este es un resumen y mis apuntes de estudio para el curso XS0100-Fundamentos de teoría estadística. Por favor no los sustituya con el libro de texto o las clases sincrónicas del curso; pues esto es un material complementario de la materia del curso. Siendo esta su última versión compilada [\[1 de julio de 2021\]](#)

Clase 1: Estimación Puntual

Fecha: 17/05/2021

Profe: Alexander Franck

Por: Ricardo Huapaya

Partimos de lo que conocimos y probamos sobre estimadores de los parámetros poblacionales, los cuales nos dan una estimación puntual, es decir, no hemos incorporado una noción probabilística a la hora de aproximarnos al verdadero valor poblacional.

Queremos encontrar un intervalo que tenga dos condiciones:

1. Que contenga al parámetro θ
2. Su amplitud sea relativamente pequeña.

Los límites de ese intervalo los llamaremos límites de confianza.

Un intervalo de confianza contiene en sus límites (variables aleatorias) el parámetro fijo (θ) con una probabilidad que se denota con $1 - \alpha$. Esta probabilidad, $1 - \alpha$ indicará la proporción de veces que la estimación de una muestra aleatoria caerá en ese intervalo.

$$P[\theta_L \leq \theta \leq \theta_U] = 1 - \alpha \quad (1.1)$$

El intervalo $[\theta_L, \theta_U]$ es el intervalo de confianza bilateral.

También podemos definir el intervalo de confianza unilateral de la forma

$$P[\theta_L \leq \theta] = 1 - \alpha \quad (1.2)$$

En este caso el intervalo está formado por un solo límite, pero tendría la forma: $[\theta_L, \infty[$. Análogamente podríamos obtener:

$$P[\theta_U \geq \theta] = 1 - \alpha$$

Con intervalo de confianza $[-\infty, \theta_u[$

1.1. Métodos

1.1.1. Método Pivotal

El método que funciona conceptualmente de base para todas las estimaciones por intervalos es el método pivotal. Este consiste en *encontrar* una variable pivote que cumpla con las siguientes dos condiciones

- Que sea función de los valores muestrales y el parámetro desconocido (θ), y este es el único parámetro desconocido.

- Que su distribución de probabilidad no dependa del parámetro (θ).

Luego partimos de las siguientes dos propiedades algebraicas:

Hecho 1.1.1. Sea c una constante arbitraria desconocida, $c > 0$, y $P(a \leq Y \leq b) = 1 - \alpha$ un intervalo de confianza para Y ; entonces note que:

1. $P(ca \leq cY \leq cb) = 1 - \alpha$
2. $P(c + a \leq c + Y \leq c + b) = 1 - \alpha$

Ejemplo 1.1.2 (8.4 del Mendelhall). Suponga que obtenemos una sola observación Y de una distribución exponencial con media θ . Use Y para construir un intervalo de confianza para θ con un coeficiente de confianza de .90.

La función de densidad de probabilidad esta dada por:

$$f(y) = \begin{cases} \left(\frac{1}{\theta}\right) e^{-\frac{y}{\theta}}, & y \geq 0 \\ 0, & y < 0 \end{cases}$$

Para ello tome $U = \frac{Y}{\theta}$, entonces:

$$f_U(u) = \begin{cases} e^{-u}, & u \geq 0 \\ 0, & u < 0 \end{cases}$$

La distribución de U no depende de θ . Entonces, podemos emplear $U = \frac{Y}{\theta}$ como cantidad pivote. Como buscamos un estimador de intervalo con coeficiente de confianza igual a .90, encontramos dos números a y b tales que:

$$P(a \leq U \leq b) = 0,90$$

Para ello debe elegir un $a \wedge b$ de la forma:

$$\begin{aligned} P(U < a) &= \int_{-\infty}^a e^{-u} du = 0,05 & P(U > b) &= \int_b^{\infty} e^{-u} du = 0,05 \\ 1 - e^{-a} &= 0,05 & e^{-b} &= 0,05 \\ a &= 0,051 & b &= 2,996 \\ \therefore & & & \end{aligned}$$

$$\begin{aligned} P(0,051 \leq U \leq 2,996) &= 0,90 = P\left(0,051 \leq \frac{Y}{\theta} \leq 2,996\right) \\ P\left(\frac{0,051}{Y} \leq \frac{1}{\theta} \leq \frac{2,996}{Y}\right) &= 0,90 \\ P\left(\frac{Y}{0,051} \geq \theta \geq \frac{Y}{2,996}\right) &= 0,90 \\ P\left(\frac{Y}{2,996} \leq \theta \leq \frac{Y}{0,051}\right) &= 0,90 \end{aligned}$$

Entonces, vemos que $Y/2,996$ y $Y/0,051$ forman los límites de confianza inferior y superior, respectivamente, que estábamos buscando. Para obtener los valores numéricos de estos límites debemos observar un valor real para Y y sustituirlo en las fórmulas dadas para los límites de confianza. Sabemos que límites de la forma $Y/2,996, Y/0,051$ incluirán los valores (desconocidos) verdaderos de θ para 90 % de los valores de Y que obtendríamos por muestreo repetido a partir de esta distribución exponencial■

Clase 2: Intervalos para medias y varianzas poblacionales

Fecha: 20/05/2021

Profe: Alexander Franck

Por: Ricardo Huapaya

2.1. Métodos

2.1.1. Métodos para la Media

Considere a X es una normal con σ^2 conocida, usamos \bar{x} para estimar μ .

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$P\left(\mu - 1,96 \frac{\sigma}{\sqrt{n}} \leq \bar{x} \leq \mu + 1,96 \frac{\sigma}{\sqrt{n}}\right) = 0,95$$

$$P\left(\bar{x} - 1,96 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + 1,96 \frac{\sigma}{\sqrt{n}}\right) = 0,95$$

Considere el caso de una σ^2 desconocida.

$$P\left(\bar{x} - t_{\alpha/2, n-1} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{x} + t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}\right) = 1 - \alpha \quad (2.1)$$

$$\text{con } t_{n-1} = \frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \quad (2.2)$$

2.1.2. Métodos para la Varianza

Defina: $Q = \frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2$

$$P\left(q_1 \leq \frac{(n-1)s^2}{\sigma^2} \leq q_2\right) = 1 - \alpha \quad (2.3)$$

$$P\left(\frac{(n-1)s^2}{q_2} \leq \sigma^2 \leq \frac{(n-1)s^2}{q_1}\right) = 1 - \alpha \quad (2.4)$$

Note que para los caso unilaterales es de la forma:

$$P(Q \geq q_1) = \frac{\alpha}{2} \wedge P(Q \leq q_2) = \frac{\alpha}{2}$$

Ejemplo 2.1.1. Un fabricante afirma que tiene un artefacto que produce mediciones con una exactitud de 1.5 micras. Con base en esta afirmación, el fabricante otorga la garantía. Se obtuvo una muestra de cinco lecturas de un mismo objeto y cuyas mediciones son las siguientes: 135, 137, 138, 137 y 139. Determine un intervalo de confianza de 90 % para la varianza poblacional.

Solución:(Excel) Note que tenemos una muestra con un tamaño de 5 observaciones, para ello podemos calcular la varianza tabulando los datos en excel para ello la formula VAR () y encontramos que s^2 de la muestra es 2,2.

Con ello nada más queda encontrar los quintiles con la fórmula de excel INV.CHICUAD (p, k_g.l.) que devuelve el inverso de la probabilidad de cola izquierda de la distribución chi cuadrado. Así $q_1 = 0,71 \wedge q_2 = 9,48$.

Finalizamos sustituyendo:

$$P\left(\frac{(5-1)2,2}{9,48} \leq \sigma^2 \leq \frac{(5-1)2,2}{0,71}\right) = 0,90$$

$$P(0,93 \leq \sigma^2 \leq 12,38) = 0,90$$

Por lo tanto para cada 100 muestras que se efectúen el 90 de estas su varianza muestral se encuentran dentro de [0.93, 12.38]. ■

2.1.3. Métodos para comparación de medias

Sea X_1, X_2, \dots, X_m muestra aleatoria de tamaño m con distribución normal con media μ_x varianza σ^2 , y sea Y_1, Y_2, \dots, Y_n muestra aleatoria de tamaño n con distribución normal con media μ_y y varianza σ^2 .

$$Q = \frac{\bar{x} - \bar{y}}{s_p \sqrt{\frac{1}{m} + \frac{1}{n}}} = \frac{\bar{x} - \bar{y}}{S_{\bar{x}-\bar{y}}}$$

Con: $s_p^2 = \frac{(m-1)s_x^2 + (n-1)s_y^2}{m+n-2}$

Por ello tenemos entonces:

$$P(-t_{\alpha/2}m + n - 2 \leq Q \leq t_{\alpha/2}m + n - 2) = 1 - \alpha \quad (2.5)$$

$$P((\bar{x} - \bar{y}) - t_{\alpha/2, m+n-2} S_{\bar{x}-\bar{y}} \leq \bar{x} - \bar{y} \leq (\bar{x} - \bar{y}) + t_{\alpha/2, m+n-2} S_{\bar{x}-\bar{y}}) = 1 - \alpha \quad (2.6)$$

2.1.4. Método para la Comparación de Varianzas

$$F = \frac{\frac{s_1^2}{\sigma_1^2}}{\frac{s_2^2}{\sigma_2^2}} \sim F_{n_1-1, n_2-1} \quad (2.7)$$

$$P\left(\frac{s_1^2/s_2^2}{F_{1-\alpha/2}} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{s_1^2/s_2^2}{F_{\alpha/2}}\right) \quad (2.8)$$

2.1.5. Métodos para una proporción y diferencia de proporciones

$$P\left(\hat{p} - z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) \quad (2.9)$$

$$P\left((\hat{p}_x - \hat{p}_y) - z_{\alpha/2}\sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{m} + \frac{\hat{p}_y(1-\hat{p}_y)}{n}} \leq \hat{p}_x - \hat{p}_y \leq (\hat{p}_x - \hat{p}_y) + z_{\alpha/2}\sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{m} + \frac{\hat{p}_y(1-\hat{p}_y)}{n}}\right) \quad (2.10)$$

2.2. Practica

Ejercicio 2.2.1. (8.50 Mendenhall) Consulte el Ejemplo 8.8. En este ejemplo, p_1 y p_2 se usaron para denotar las proporciones de refrigeradores de las marcas A y B, respectivamente, que fallaron durante los períodos de garantía

1. En el nivel aproximado de 98 % de confianza, ¿cuál es el mayor “valor creíble” para la diferencia en las proporciones de fallas de refrigeradores de las marcas A y B?
2. En el nivel aproximado de 98 % de confianza, ¿cuál es el menor “valor creíble” para la diferencia en las proporciones de fallas de refrigeradores de las marcas A y B?
3. Si $p_1 - p_2$ es realmente igual a 0.2251, ¿cuál marca tiene la mayor proporción de fallas durante el período de garantía? ¿Qué tanto más grande?
4. Si $p_1 - p_2$ es realmente igual a -0.1451, ¿cuál marca tiene la mayor proporción de fallas durante el período de garantía? ¿Qué tanto más grande?
5. Como se observó en el Ejemplo 8.8, cero es un valor creíble de la diferencia. ¿Concluiría usted que hay evidencia de una diferencia en las proporciones de fallas (dentro del período de garantía) para las dos marcas de refrigeradores? ¿Por qué?

Solución

Podemos reescribir el Intervalo de Confianza para las proporciones de la forma,

$$(\hat{p}_x - \hat{p}_y) \pm z_{\alpha/2}\sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{m} + \frac{\hat{p}_y(1-\hat{p}_y)}{n}}$$

Usando excel podemos hacer el calculo de $Z_{\alpha/2}$ INV.NORM.ESTAND(0,99), en este caso nos da 2,33 con ello podemos calcular el coeficiente de la forma.

$$(0,24 - 0,20) \pm 2,33\sqrt{\frac{0,24(1-0,24)}{50} + \frac{0,20(1-0,20)}{60}}$$

Por lo que el intervalo de confianza nos da $[-0,145, 0,225]$.

(1) En el nivel 98 % de confianza el mayor valor creíble es 0,225. (2) De forma análoga el menor valor creíble de diferencia entre las proporciones es de -0.145.

(3) En este caso la marca A sería la que presente más fallos y sería casi el doble de fallos (4) dado caso la marca B sería la que más presente fallos y sería cuatro veces la cantidad de fallo. (5) No hay evidencia pues 0 es un valor dentro del intervalo.

Ejercicio 2.2.2. 33

Ejercicio 2.2.3. (8.58 Mendenhall) Los administradores de un hospital deseaban estimar el número promedio de días necesarios para el tratamiento de enfermos internados entre las edades de 25 y 34 años. Una muestra aleatoria de 500 pacientes entre estas edades produjo una media y una desviación estándar igual a 5.4 y 3.1 días, respectivamente. Construya un intervalo de confianza del 95 % para la duración media de permanencia de la población de pacientes de la cual se extrajo la muestra.

Solución: (Excel)

Para ello aplicamos

$$P\left(\bar{x} - 1,96 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + 1,96 \frac{\sigma}{\sqrt{n}}\right) = 0,95$$

Sustituimos valores:

$$P\left(5,4 - 1,96 \frac{3,1}{\sqrt{500}} \leq \mu \leq 5,4 + 1,96 \frac{3,1}{\sqrt{500}}\right) = 0,95$$

$$P(5,13 \leq \mu \leq 5,67) = 0,95$$

En este caso el valor de 1,96 es el $z_{\alpha/2}$ que se obtiene en excel aplicando la formula:

INV.NORM.ESTAND(I.de Confianza + (1-I.de Confianza)/2)

Ejercicio 2.2.4. (8.62 Mendenhall) Los siguientes estadísticos son el resultado de un experimento realizado por P. I. Ward para investigar una teoría relativa al comportamiento de cambio de piel del macho *Gammarus pulex*, un pequeño crustáceo. Si el macho cambia de piel mientras se aparea con una hembra, éste debe liberarla y perderla. La teoría es que el macho *Gammarus pulex* es capaz de posponer dicho cambio, con lo cual reduce la posibilidad de perder su pareja. Ward asignó aleatoriamente 100 parejas de machos y hembras a dos grupos de 50 cada uno. Las parejas del primer grupo se mantuvieron juntas (normal); las del segundo grupo fueron separadas. Se registró el tiempo de muda para machos y hembras, y las medias, desviaciones estándar y tamaños muestrales se ilustran en la tabla siguiente. (El número de crustáceos en cada una de las cuatro muestras es menor que 50 porque algunos en cada grupo no sobrevivieron hasta el tiempo de muda.)

Clase 3: Contraste Hipótesis estadísticas

Fecha: 20/05/2021

Profe: Alexander Franck

Por: Ricardo Huapaya

3.1. Preliminares

El objetivo de la estadística es hacer inferencias acerca de parámetros poblacionales desconocidos con base en información contenida en datos muestrales. Las inferencias se pueden interpretar como **estimaciones de parámetros** o también como las **pruebas de hipótesis de los parámetros**.

Esto nos lleva a las pruebas estadísticas lo cual el objetivo es probar una hipótesis concerniente a los valores de uno o más parámetros poblacionales. En cierto sentido hay una prueba **por contradicción**.

Elementos de una prueba estadística

1. Hipótesis nula, H_0
2. Hipótesis alternativa, H_a
3. Estadístico de prueba
4. Región de rechazo

3.1.1. Hipótesis estadísticas

Por ejemplo, si se quiere probar que $\mu = \mu_0$ son base en una muestra aleatoria X_1, X_2, \dots, X_m la **hipótesis nula** es $H_0 : \mu = \mu_0$

Esta hipótesis se contrasta con la **hipótesis alternativa** H_a que especifica un valor alternativo para μ puede ser $H_a : \mu \neq \mu_0$ o $H_a : \mu > \mu_0$ o $H_a : \mu < \mu_0$

Hipótesis simple: se especifica completamente la distribución de la cual se tomó la muestra, $H_0 : \mu = \mu_0$

Hipótesis compuesta: no se especifica completamente la distribución de la cual se tomó la muestra, por ejemplo, $H_a : \mu > \mu_0$

3.1.2. Tipos de Error

Cuando contrastamos H_0 vs H_1 , podemos cometer dos tipos de error:

1. Rechazar H_0 cuando esta es verdadera: Error tipo I $P(I) = \alpha$. Se llama nivel de significancia del contraste.
2. No rechazar H_0 cuando esta es falsa: Error tipo II $P(II) = \beta$

3.1.3. Estadístico y Región de Rechazo

Necesitamos un estimador de θ , $\hat{\theta} = h(x_1, x_2, \dots, x_n)$ y determinar su distribución muestral $g(\hat{\theta}; \theta_0)$.

Dividir la región de todos los valores posibles \mathcal{R} en dos regiones: una región de rechazo RR , en la que θ_0 es poco probable, tal que $P(\hat{\theta} \in RR | \theta = \theta_0) = P(I) = \alpha$.

Calcular el valor de $\hat{\theta}$ con la muestra aleatoria y rechazar H_0 si cae en la región de rechazo, y, si no, no rechazar H_0 .

Ejemplo 3.1.1. (Mendelhall 10.1) Para la encuesta política de Jones se muestrearon $n = 15$ votantes. Deseamos probar $H_0 : p = ,5$ contra la alternativa, $H_a : p < ,5$. El estadístico de prueba es Y , el número de votantes muestreados a favor de Jones. Calcule a si seleccionamos región de rechazo $RR = \{y \leq 2\}$ como la región de rechazo.

Solucion: Considere que Y es una var aleatoria binomial, pues votan por candidato o no. Entonces considere:

$$\begin{aligned}\alpha &= \text{Probabilidad de cometer el error tipo I} = \text{Probabilidad de rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera} \\ &= P(\text{valor del estadístico de prueba esta en } RR \text{ cuando } H_0 \text{ es verdadera}) \\ &= P(Y \leq 2 \text{ cuando } p = 0,5)\end{aligned}$$

Asi bien con $n = 15$, $p = 0,50$ tenemos la distribución acumulada de la forma

$$\begin{aligned}\alpha &= P(Y \leq 2) = \sum_{y=0}^2 \binom{15}{y} (.5)^y (.5)^{1-y} \\ &= 0,004\end{aligned}$$

En conclusión si se usa la región de rechazo $RR = \{y \leq 2\}$ asumimos que el riesgo de rechazar la hipótesis nula es muy bajo.

Ejercicio 3.1.2. Un profesor universitario quiere conocer una nueva forma de enseñar la materia. Anteriormente solo el 0,4 de los estudiantes pasaba el curso, por lo que interesa contrastar la hipótesis:

$$H_0 : p = 0,4 \wedge H_1 : p > 0,4$$

Para ellos, el profesor toma una muestra de 12 estudiantes y observa cuanto de ellos aprueban en el curso.

Un colega del profesor le dice que el puede pude tener el mismo nivel de significancia al rechazar la hipótesis nula cuando más de tres estudiantes pasan el curso y si al tirar una moneda cuatro veces se observa escudo en todos los lanzamientos. Calcule el valor de significancia.

Solución:

Considere el estadístico Y como una binomial que en este caso es de la forma $Y \sim \text{Bin}(12, 0,4)$. Entonces α nivel de significancia es de la forma:

$$\begin{aligned}\alpha &= P(Y > 3) \cdot P(\text{Cuatro escudos}) \\ &= (1 - P(Y \leq 3)) \cdot P(\text{Cuatro escudos}) \\ &= \sum_{y=0}^3 \binom{12}{y} (.4)^y (.6)^{12-y} \cdot (.5)^4 \\ &= 0,775 * 0,062 = 0,048\end{aligned}$$

Para encontrar a la distribución de la binomial se aplica la siguiente formula de excel:

$$1-DISTR.BINOM.N(3,12,0.4,VERDADERO)$$



3.2. Procedimientos de contraste de hipótesis

3.2.1. Pruebas comunes con muestras grandes

Asuma que se quiere probar un conjunto de hipótesis con respecto a un parámetro θ con en una muestra aleatoria de $listnY$. Para desarrollar los procedimiento de prueba se debe basar en un estimador $\hat{\theta}$ con distribución normal con media θ y error estándar $\sigma_{\hat{\theta}}$.

Si se tiene un θ_0 es un valor posible de θ , podemos plantear la prueba con la hipótesis nula de $H_0 : \theta = \theta_0$ y la hipótesis alternativa $H_a : \theta > \theta_0$ con una región de rechazo que definiremos de la forma $RR = \{\hat{\theta} > k\}$. Donde k se determina al fijar la probabilidad de que se cometa el error de tipo I.

Considere entonces:

$$k = \theta_0 + z_\alpha \sigma_{\hat{\theta}}$$

Siendo k la selección apropiada si Z tiene una distribución normal estándar, entonces z_α es tal que $P(Z > z_{\alpha}) = \alpha$. Así pues:

$$RR = \{\hat{\theta} : \hat{\theta} > \theta_0 + z_\alpha \sigma_{\hat{\theta}}\} = \left\{ \hat{\theta} : \frac{\hat{\theta} - \theta_0}{\sigma} \right\}$$

Entonces tome como estadístico de prueba $Z = \frac{\hat{\theta} - \theta_0}{\sigma}$ y así podemos reescribir la región de rechazo de la forma: $RR = \{z > z_\alpha\}$

En este caso estamos hablando de un región de **rechazo en la cola superior**; note que este tipo se utiliza en su mayoría cuando nos interesa detectar un aumento.

Pruebas de Hipótesis de nivel α para muestras grandes. ($n \geq 30$) :

$$H_0 : \theta = \theta_0 \tag{3.1}$$

$$H_a : \begin{cases} \theta > \theta_0 & \text{cola superior} \\ \theta < \theta_0 & \text{cola inferior} \\ \theta \neq \theta_0 & \text{ambas colas} \end{cases} \tag{3.2}$$

$$\text{Estadístico de prueba } Z = \frac{\hat{\theta} - \theta_0}{\sigma} \tag{3.3}$$

$$RR : \begin{cases} z > z_\alpha & \text{cola superior} \\ z < -z_\alpha & \text{cola inferior} \\ |z| > z_{\alpha/2} & \text{ambas colas} \end{cases} \tag{3.4}$$

Ejemplo 3.2.1. (Media) El vicepresidente de ventas de una gran empresa afirma que los vendedores están promediando no más de 15 contactos de venta por semana. (Le gustaría aumentar esta cantidad.) Como prueba de su afirmación, aleatoriamente se seleccionan $n = 36$ vendedores y se registra el número de contactos hechos por cada uno para una sola semana seleccionada al azar. La media y varianza de las 36 mediciones fueron

17 y 9, respectivamente. ¿La evidencia contradice lo dicho por el vicepresidente? Use una prueba con nivel $\alpha = ,05$.

Solución: Aplicamos una prueba de hipótesis con cola superior pues a vicepresidente le interesa rechazar la media más pequeña ya que la grande creció. Aplicamos la formula para poner a prueba lo siguiente.

$$H_0 : \mu_0 = 15 \wedge H_1 : \mu > 15$$

Entonces, estableciendo la región de rechazo con $\alpha = 0,05$. Esta dada por $\{z > z_\alpha\} = \{z > z_{0,05}\} = \{z > 1,64\}$. Vemos que,

$$\begin{aligned} Z &= \frac{\hat{\theta} - \theta_0}{\sigma} \\ &= \frac{\mu - \mu_0}{\frac{\sigma}{\sqrt{n}}} \\ &= \frac{17 - 15}{\frac{3}{\sqrt{36}}} \\ &= 4 \end{aligned}$$

Por lo que concluimos que, al nivel de significancia $\alpha = ,05$., la evidencia es suficiente para indicar que la afirmación del vicepresidente es incorrecta y que el número promedio de contactos de ventas por semana es mayor que 15. Se rechaza la hipótesis nula. ■

Pruebas de Hipótesis de nivel α para para proporciones :

$$H_0 : p = p_0 \quad (3.5)$$

$$H_a : \begin{cases} p > p_0 & \text{cola superior} \\ p < p_0 & \text{cola inferior} \\ p \neq p_0 & \text{ambas colas} \end{cases} \quad (3.6)$$

$$\text{Estadístico de prueba } Z = \frac{p - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \quad (3.7)$$

$$\text{RR} : \begin{cases} z > z_\alpha & \text{cola superior} \\ z < -z_\alpha & \text{cola inferior} \\ |z| > z_{\alpha/2} & \text{ambas colas} \end{cases} \quad (3.8)$$

Pruebas de Hipótesis de nivel α para para proporciones :

$$\text{Estadístico de prueba } Z = \frac{(p_{e1} - p_{e2}) - (p_1 - p_2)}{\sqrt{p_e(1-p_e)}\sqrt{\frac{1}{n} + \frac{1}{m}}} \quad (3.9)$$

$$p_e = \frac{np_{e1} + mp_{e2}}{m + n} \quad (3.10)$$

$$\text{RR} : \begin{cases} z > z_\alpha & \text{cola superior} \\ z < -z_\alpha & \text{cola inferior} \\ |z| > z_{\alpha/2} & \text{ambas colas} \end{cases} \quad (3.11)$$

3.2.2. Pruebas para muestras pequeñas

3.2.3. Niveles de significancia y p-values

Una vez tomada una decisión sobre un estadístico de prueba, a veces es posible presentar el valor p o el nivel de significancia alcanzado y que está relacionado con una prueba. Esta cantidad es un estadístico que representa el valor más pequeño de α para el cual se puede rechazar la hipótesis nula.

Definition 3.2.2. (P-Value) Sea W es un estadístico de prueba, el valor p , o nivel de significancia alcanzado, es el nivel más pequeño de significancia α para el cual la información observada indica que la hipótesis nula debe ser rechazada.

El método general para calcular valores p . Si fuéramos a rechazar H_0 en favor de H_a para valores pequeños de un estadístico de prueba W , por ejemplo.

$$RR : \{w \leq k\}$$

El valor p relacionado con un valor observado w_0 de W está dado por:

$$\text{valor } p = P(W \leq w_0, \text{ cuando } H_0 \text{ es verdadera})$$

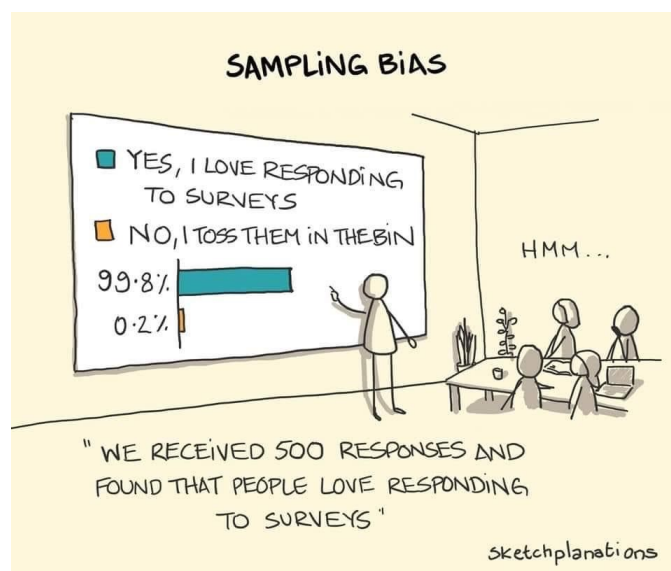


Figura 3.1: Un meme para no olvidar que estadística es lindo

Clase 4: Lemma Neyman Pearson

Fecha: 21/05/2021

Profe: Alexander Franck

Por: Ricardo Huapaya

4.1. Concepto de Espacio paramétrico y potencia

Sea una población descrita por $f_{x;\theta}$ y Ω el conjunto de valores que puede tomar θ , llamado espacio paramétrico. Sea ω el subconjunto de Ω que contiene todos los valores de θ que establece la hipótesis nula.

$$\begin{aligned}H_0 : \theta &= \theta_0 \\H_a : \theta &\in \Omega - \{\theta_0\}\end{aligned}$$

Sea $\hat{\theta}$ un estimador de θ que se utiliza para contrastar H_0 con H_a y RR la zona de rechazo. La función potencia $\Pi(\theta)$ definida para $\theta \in \Omega$ por:

$$\Pi(\theta) = \begin{cases} P(\hat{\theta} \in RR) = 1 - \beta & \text{si } \theta \neq \theta_0 \\ \alpha & \text{si } \theta = \theta_0 \end{cases} \quad (4.1)$$

4.2. Neyman-Pearson

Con un nivel de α fijado, queremos maximizar la potencia de la prueba. Para eso se usa el lema de Neyman-Pearson (solo para contrastar hipótesis simples)

Lemma 4.2.1. (Neyman-Pearson) Para contrasta $H_0 : \theta = \theta_0$ con $H_a : \theta = \theta_a$, con base en una muestra aleatoria Y_1, Y_2, \dots, Y_n de una distribución con parámetro θ . Sea $L(\theta)$ la función de verosimilitud de la muestra cuando el valor del parámetro es θ . Entonces, para un α dado, la prueba que maximiza la potencia en θ_a tiene una región de rechazo RR determinada por:

$$\frac{L(\theta_0)}{L(\theta_a)} < k$$

El valor de k se escoge para que tenga el valor deseado para α

Clase 5: Estadística Bayesiana

Fecha: 28/6/2021

Profe: Alexander Franck

Por: Ricardo Huapaya

5.1. Contraste con el Enfoque frecuentista

Coincide con los enfoques de la probabilidad, incluido los enfoques **a priori** racionalista de inferencias y estocástico determinado por la cantidad de resultados; **el enfoque clásico** sea frecuentista que de manera escéptica el enfoque de inferencia y estima a partir del estudio empírico usando las muestras.

Enfoque subjetivo: basado en creencias dado a que incorpora información subjetiva en el cálculo matemático, considera que θ debe ser una distribución de probabilidad. Pues así antes de ver la probabilidad de los datos asignamos una distribución a priori que llamamos $\Pi(\theta)$, $0 \leq \theta \leq 1$. Queremos encontrar una distribución a posteriori.

La distribución posteriori $\Pi(\theta|y)$ se obtiene por medio del teorema de Bayes.

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B)} \quad (5.1)$$

$$\Pi(\theta|Y) = \frac{P(Y|\theta) \cdot P(\theta)}{P(Y)} \quad (5.2)$$

También puedo escribir la función de distribución posterior como:

$$\Pi(\theta|y) \propto f(y|\theta) \cdot \Pi(\theta)$$

Ejemplo 5.1.1. Considere la función de verosimilitud de un solo lanzamiento de dados se puede escribir como $P(Y_i|\theta) = \theta^{y_i}(1-\theta)^{1-y_i}$, es decir $P(Y_i = 1) = \theta$ y $P(Y_i = 0) = 1 - \theta$.

Entonces,

$$\begin{aligned} P(Y_1|\theta, \dots, Y_n|\theta) &= \theta^{y_1}(1-\theta)^{1-y_1} \dots \theta^{y_n}(1-\theta)^{1-y_n} \\ &= \prod_{i=1}^n \theta^{y_i}(1-\theta)^{1-y_i} \\ &= \theta^{\sum y_i} (1-\theta)^{n-\sum y_i} \end{aligned}$$

La especificación del modelo requiere una distribución a priori de $0 \leq \theta \leq 1$. ■

