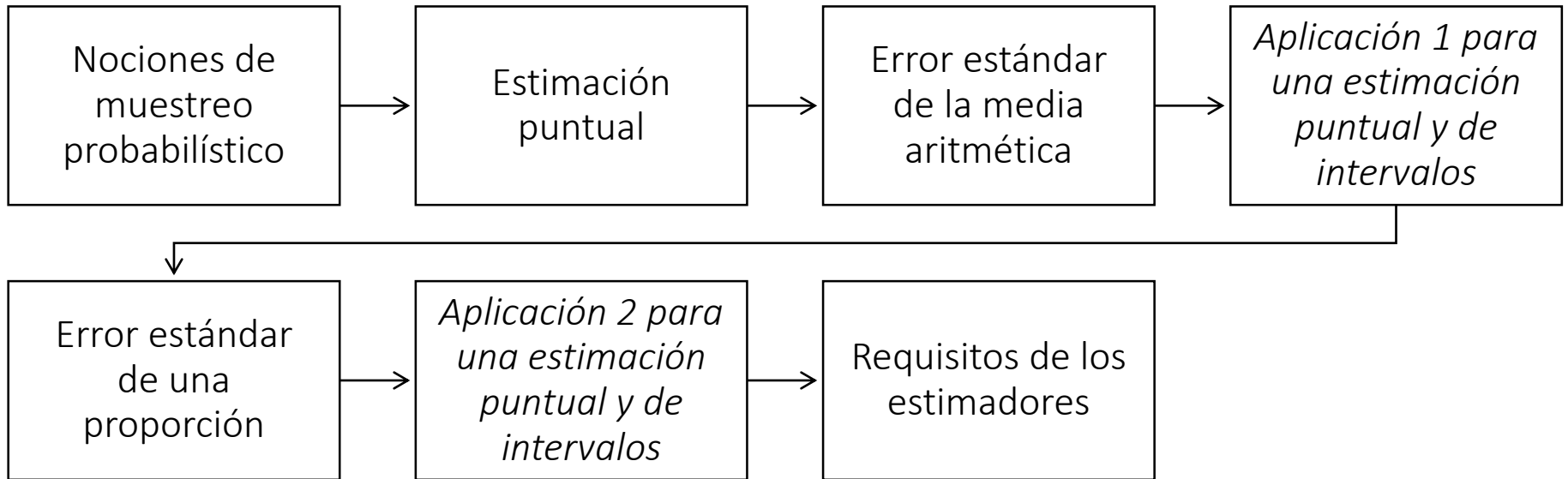


Distribuciones Muestrales e Intervalos de Confianza

Sesión 5

Lo que vamos a ver hoy



Muestras y representatividad

- La mayor parte de la información que poseemos proviene de muestras, algunas probabilísticas y otras no
- La mayor parte de las decisiones que tomamos las realizamos sobre muestras, es decir, con M y s
- Pregunta siempre presente: ¿es la muestra adecuada para conocer el parámetro?
- Contestar estas 3 preguntas:
 - ¿La muestra corresponde a la población bajo estudio?
 - ¿Qué método se utilizó para seleccionar los casos (personas) de esta muestra?
 - ¿Son los casos representativos de la población?
 - Un estudio que no permite contestar estas preguntas es poco confiable
- Sobre la representatividad: significa que la muestra representa fielmente las características (variabilidad) del universo
 - Muestra representativa: la variabilidad entre sus miembros es similar a la variabilidad de todos los miembros del universo

Muestreo probabilístico

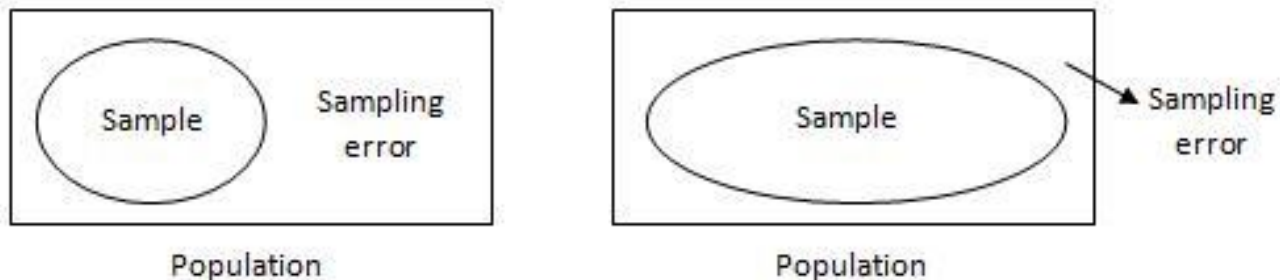
- Lógica:
 - Se basa en un proceso de selección aleatorio
 - Todos tienen una probabilidad conocida de ser seleccionados
 - Ventaja principal: Alta capacidad de inferencia
 - Invariablemente, hay error muestral y hay riesgo de que la muestra no sea representativa; esto por simple azar
- Principios:
 1. Cuanto más grande la muestra, más confianza podemos tener de que la muestra es representativa
 2. Cuanto más homogénea sea la población bajo estudio, más confianza podemos tener de que la muestra es representativa; esto sea del tamaño que sea.

Población diversa → muestra grande

Población similar → muestra pequeña

Muestreo probabilístico

- El error muestral es la diferencia entre el estadístico y el parámetro
 - Ejemplo: $M - \mu$
- Cómo reducirlo?
 - Asegurando la representatividad
 - Aumentando el tamaño de muestra ($n \rightarrow N$)
- Ley de los Grandes Números:
 - La media muestral (M) se acerca a media poblacional (μ) conforme aumenta “ n ” y se acerca a “ N ”
 - El tamaño de la muestra se acerca al tamaño de la Población
 - Por ende la fuerza de una inferencia depende de la cantidad de información disponible y variación real capturada/medida



Tipos de poblaciones y muestreos

- Estadísticamente, las poblaciones son de 2 tipos:
 - Finitas: el tamaño es conocido y se posee el marco muestral
 - Aplica un muestreo aleatorio simple (SRS) equi-probable (prob. clásica)
 - Esto porque se puede saber el número total de muestras disponibles

$$\frac{N!}{n!(N - n)!}$$

- Aquí aplica la fórmula de cálculo de tamaño de muestra simple
- Infinitas: el tamaño es desconocido (no se posee el marco muestral) o muy grande
 - Igualmente hay que asegurar la representatividad evitando el sesgo de selección
 - Aplica la fórmula de cálculo de tamaño de muestra compleja

Estimación puntual

- Los parámetros son desconocidos y los estadísticos conocidos
- Los estimadores son los estadísticos muestrales
- El muestreo nos sirve para:
 - Obtener estadísticos (estimadores puntuales) e intervalos de confianza (intervalo estimado o límites probabilísticos dentro del cual el parámetro puede encontrarse)
 - Realizamos estimaciones de parámetros puntuales y por intervalos
 - Los intervalos agregan información sobre el posible error de estimación
- Cálculo: la media (M) y la desviación estándar (s) de la muestra son los estimadores puntuales de los parámetros respectivos de la población

$$M = \mu \quad s = \sigma \quad \bar{p} = \pi$$

- Error estándar: mide la precisión del estimador puntual con relación al parámetro
 - Menor error estándar → más preciso es el estimador
 - Mayor error estándar → menos preciso es el estimador

El error estándar de la media aritmética (M)

- Cálculo:

- Utilizamos los datos que nos da la muestra
- Difiere según la población sea finita o infinita:

$$\textit{Finita: } s_M = \sqrt{\frac{N-n}{N-1}} * \left(\frac{s}{\sqrt{n}} \right) \quad \textit{Infinita: } s_M = \left(\frac{s}{\sqrt{n}} \right)$$

- Lógica:

- Mientras más grande sea la muestra (n) menor será el error estándar (el estimador puntual será más preciso en relación con el parámetro)
- Mientras más homogénea o similar sea la muestra (o población) menor será el error estándar (mayor precisión)

Utilidad en Políticas Públicas

- Vilalta, C. 2007. Compra y coacción del voto en México: variaciones estatales y diferencias de opinión
 - Liga: http://www.programassociales.org.mx/biblioteca/Serie_ENAPP_No_4_Vilalta.pdf
- Pregunta: ¿Cuál es la cantidad estimada en pesos con que se compra el voto en México?
 - Fuente: Encuesta Nacional Sobre la Protección de Programas Sociales Federales (ENAPP-2006)
 - Nota: el reactivo a esa pregunta es “¿Por cuánto dinero piensa usted que otras personas intercambiarían su voto? Esta es una medida subjetiva e indirecta
- Veamos los estimadores puntuales y el error estándar (precisión) de esa estimación

Utilidad en Políticas Públicas

- Los resultados de la encuesta (en abril de 2006) fueron los siguientes:
 - Media aritmética (M): 88.2 pesos
 - Desviación estándar (s): 29.3 pesos
 - $n = 332$
- ¿Cuál es el error estándar?
 - Asumiendo una población infinita:

$$s_M = \left(\frac{29.3}{\sqrt{332}} \right) = \frac{29.3}{18.2} = 1.6 \text{ pesos}$$

- Equivalente 2% (es decir 1.6/88.2) de la media muestral... ¿Es mucho o poco?
- Asumiendo una población finita:

$$s_M = \sqrt{\frac{71,730,536}{71,730,867} \left(\frac{29.3}{18.2} \right)} = 1.0 * 1.6 = 1.6 \text{ pesos}$$

Utilidad en Políticas Públicas

- Para responder si es mucho o poco, es preferible realizar una estimación por intervalo de confianza (IC)
- Este se calcula de la siguiente manera:

$$M - 1.96 \left(\frac{s}{\sqrt{n}} \right) \leq \mu \leq M + 1.96 \left(\frac{s}{\sqrt{n}} \right)$$

The diagram illustrates the components of the confidence interval formula. The left side of the inequality, $M - 1.96 \left(\frac{s}{\sqrt{n}} \right)$, is broken down into M (labeled NC) and $1.96 \left(\frac{s}{\sqrt{n}} \right)$ (labeled ES). The right side, $M + 1.96 \left(\frac{s}{\sqrt{n}} \right)$, is broken down into M (labeled NC) and $1.96 \left(\frac{s}{\sqrt{n}} \right)$ (labeled ES). Brackets below these terms group them into ME (Marginal Error) for both sides. A large bracket at the bottom encompasses the entire inequality, with a box below it containing the text "INTERVALO DE CONFIANZA SIMÉTRICO Y NORMAL ESTÁNDAR".

INTERVALO DE CONFIANZA SIMÉTRICO Y NORMAL ESTÁNDAR

- Donde:
 - M: media aritmética de la muestra
 - s: desviación estándar de la muestra
 - n: tamaño de la muestra
 - 1.96: valor crítico de Z con un NC del 95% (para una distribución normal estándar)

Ejemplo de IC para una media aritmética

- Los resultados de la encuesta fueron:
 - $M = 88.2$ pesos / $s = 29.3$ pesos / $n = 332$
- Para un nivel de confianza del 95% ($NC = 1.96$), esto se calcularía de la siguiente manera:

$$M - 1.96 \left(\frac{s}{\sqrt{n}} \right) \leq \mu \leq M + 1.96 \left(\frac{s}{\sqrt{n}} \right)$$

$$88.2 - 1.96 \left(\frac{29.3}{\sqrt{332}} \right) \leq \mu \leq 88.2 + 1.96 \left(\frac{29.3}{\sqrt{332}} \right)$$

- Es decir, según la información proveniente de la muestra, la opinión es que la gente vendería su voto entre:

$$85.05 \leq \mu \leq 91.35$$

Con una probabilidad de error del 5% ($1 - NC = 1 - 0.95 = 0.05$)

Ejemplo de IC para una media aritmética

- Si deseáramos un mayor nivel de confianza, esto implicaría incrementar el IC
 - Ojo: el error estándar seguiría siendo el mismo (1.6 pesos)
- Por ejemplo, si deseáramos tener un IC de μ con un nivel de confianza del 99%, esto implicaría que el área de probabilidades o valores posibles del estimador puntual debería ser más amplio. En el caso de una curva normal, esto implicaría aumentar el intervalo de $z = 1.96$ a $Z = 2.58$

- En consecuencia:

$$88.2 - 2.58 \left(\frac{29.3}{\sqrt{332}} \right) \leq \mu \leq 88.2 + 2.58 \left(\frac{29.3s}{\sqrt{332}} \right)$$

$$84.05 \leq \mu \leq 92.35$$

- Es decir, con una probabilidad de error del 1% (NC = 99%) podemos inferir que los encuestados opinan que la gente vendería su voto entre 84 y 92 pesos

Error estándar de una proporción (p)

- Lo mismo se puede hacer para un estimador proporcional
- Cálculo:
 - Utilizamos los datos que nos da la muestra
 - También difiere según la población sea finita o infinita:

$$\textit{Finita: } s_{\bar{p}} = \sqrt{\frac{N-n}{N-1}} * \sqrt{\frac{p(1-p)}{n}} \qquad \textit{Infinita: } s_{\bar{p}} = \sqrt{\frac{p(1-p)}{n}}$$

- Aplica la misma lógica:
 - Mientras más grande sea la muestra (n) menor será el error estándar (el estimador puntual será más preciso en relación con el parámetro)
 - Mientras más homogénea o similar sea la muestra (o población) menor será el error estándar (mayor precisión)

Utilidad en Políticas Públicas

- En la misma encuesta sobre delitos electorales, hubo una proporción de encuestados beneficiarios de programas sociales que reportó haber sido coaccionado/presionado en su voto. De los 1,371 encuestados que dieron una respuesta a esa pregunta:
 - 30 ($p = .022$) reportaron sí haber sido coaccionados
 - 1,341 ($1-p = .978$) reportaron no haber sido coaccionados
- Si quisiéramos saber con un nivel de confianza del 95% ($Z = 1.96$) en qué intervalo se encuentra el porcentaje (Π) de la población que fue presionada para que votara a favor de uno u otro partido, lo calcularíamos de la siguiente manera:

$$0.022 - 1.96 \sqrt{\frac{.022(.978)}{1371}} \leq \Pi \leq 0.022 + 1.96 \sqrt{\frac{.022(.978)}{1371}}$$

$$1.42\% \leq \Pi \leq 2.98\%$$

Población que dice haber sido coaccionada en su voto

¿Entonces cuánta gente fue coaccionada?

- Si conociéramos el tamaño de la población o universo (N), podríamos aplicar el IC anterior para estimar el número de personas que pueden haber sido coaccionadas en su voto
- Pero considerando que desconocemos el número de beneficiarios podemos entonces calcular una tasa por cien mil:
 - N: 100,000
- Tendríamos que:

$$1.42\% \leq \Pi \leq 2.98\%$$

$$1,424 \leq \mu \leq 2,976$$

- Es decir, entre 1,424 y 2,976 personas fueron coaccionada de un universo de 100,000 personas; con un NC del 95%

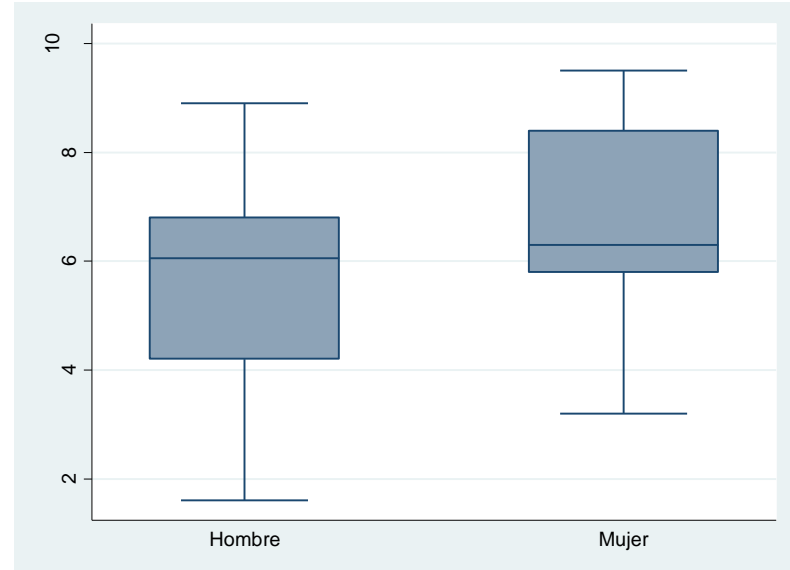
Requisitos de los estimadores

- Insegados: se espera que no haya diferencia entre el estimador puntual y el parámetro
- Eficientes: se espera que el error estándar sea pequeño
 - Comparativamente, el mejor estimador es aquel con el error estándar más pequeño
- Consistentes (robustos): se espera que conforme aumente el tamaño de la muestra, el estadístico se acerque al parámetro (no se desvíe del mismo)

Calificaciones del primer parcial

Diferencias entre Mujeres y Hombres

- ¿Quiénes salieron mejor?



Two-sample t test with unequal variances

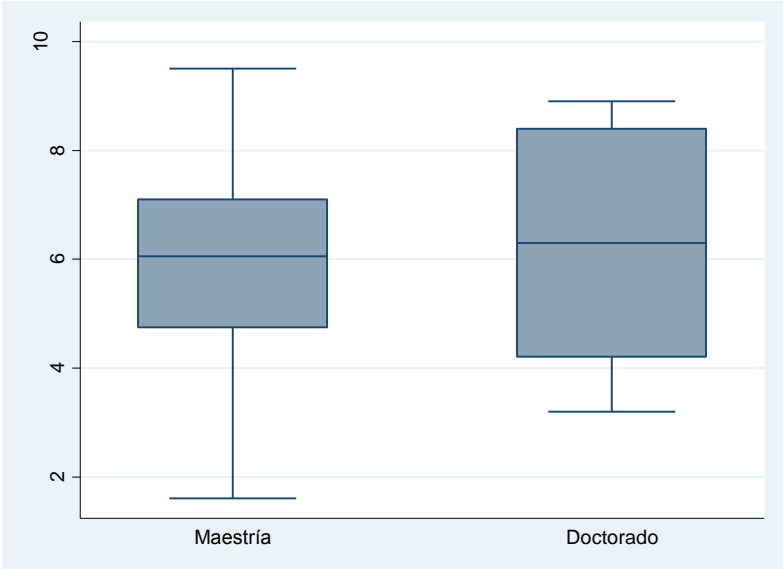
Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
Hombre	22	5.695455	.4428595	2.077195	4.774478	6.616431
Mujer	13	6.638462	.5790729	2.087877	5.37677	7.900153
combined	35	6.045714	.3552569	2.101728	5.323745	6.767683
diff		-.943007	.7290062		-2.443779	.5577651

diff = mean(**Hombre**) - mean(**Mujer**) t = **-1.2936**
Ho: diff = 0 Satterthwaite's degrees of freedom = **25.2135**

Ha: diff < 0 Ha: diff != 0 Ha: diff > 0
Pr(T < t) = **0.1038** Pr(|T| > |t|) = **0.2075** Pr(T > t) = **0.8962**

Diferencias entre Maestría y Doctorado

- ¿Quiénes salieron mejor?



Two-sample t test with equal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
Maestría	28	5.942857	.3967936	2.099635	5.128704	6.75701
Doctorad	7	6.457143	.8405942	2.224003	4.400283	8.514003
combined	35	6.045714	.3552569	2.101728	5.323745	6.767683
diff		-.5142857	.8970421		-2.339332	1.31076

diff = mean(Maestría) - mean(Doctorad) t = -0.5733
 Ho: diff = 0 degrees of freedom = 33

Ha: diff < 0 Ha: diff != 0 Ha: diff > 0
 Pr(T < t) = 0.2852 Pr(|T| > |t|) = 0.5703 Pr(T > t) = 0.7148

Tarea

- Siguiente clase: cálculo de tamaño de muestra y pruebas de hipótesis
- Anderson et al.: Resolver ejercicios (entregar a mano y con procedimientos):
 - Estimación puntual: Ejercicios 11 a 17 del Cap. 7
 - Error estándar: Ejercicios 20, 21, 30, 33 y 34 del Cap. 7
 - Intervalos de confianza: Ejercicios 1 a 10 del Cap. 8
- Entrega: jueves en el laboratorio