

Desagregación de datos: Posibilidades y limitaciones de las encuestas de hogares y métodos de estimación en áreas pequeñas

Xavier Mancero
División de Estadísticas
CEPAL



Proyecto “Enfoques innovadores para examinar la desigualdad mediante la integración de diferentes fuentes de datos en América Latina y el Caribe”



Motivación

- La implementación de los Objetivos de Desarrollo Sostenible pone énfasis en **no dejar a nadie atrás**, para superar las desigualdades que caracterizan a la región.
 - Para ello, las estadísticas deben estar desagregadas por características relevantes de la población.
 - "Los indicadores de los Objetivos de Desarrollo Sostenible deben desglosarse, cuando corresponda, por ingreso, género, edad, raza, etnia, estatus migratorio, discapacidad y ubicación geográfica; u otras características, de conformidad con los Principios Fundamentales de las Estadísticas Oficiales".
- El conocimiento sobre las desigualdades está limitado por la disponibilidad de información:
 - Encuestas de hogares subestiman la desigualdad en la distribución de los ingresos, ya que no captan adecuadamente a los hogares más ricos.
 - Capacidad limitada para proporcionar información desagregada para grupos de población y áreas geográficas específicas.
 - Aprovechamiento insuficiente de datos satelitales e información geográfica para cuantificar y visibilizar las desigualdades.

Posibilidades del uso combinado de fuentes de información

- Desarrollos metodológicos recientes para combinar información de diferentes fuentes de datos, tales como encuestas de hogares, censos de población, registros administrativos o Cuentas Nacionales.
- Uso de datos satelitales e información geográfica para producir estadísticas y representación de la información a través de mapas.
- La disponibilidad de mayor y mejor información sobre las desigualdades es un insumo valioso para las políticas públicas que apuntan a corregirlas.



Componentes del proyecto

Medición de desigualdad del ingreso con datos de EH, registros tributarios y Cuentas Nacionales

Integración de información estadística y geoespacial

Desagregación de información mediante metodologías de “estimación de áreas pequeñas”

Actividades del proyecto



Posibilidades y limitaciones
de las encuestas de hogares
para la desagregación de datos



Las encuestas de hogares como fuente de información

- Una de las principales fuentes de información sobre las condiciones de vida.
- Útil para producir desagregaciones para diferentes grupos de población.
- Una encuesta se planea para generar información para dominios de estudio predefinidos.
 - Ej.: Tasa de desempleo, por área urbana y rural
- Sin embargo, podríamos estar interesados en estimaciones para subgrupos de población que no se abordaron en el diseño.
 - Ej.: Asistencia escolar de los niños (6-12 años), por quintil de ingreso
- En estos casos, la desagregación podría no ser viable:
 - Falta de información (sin casos observados)
 - Baja precisión

Precisión de las estimaciones

- Una encuesta es una colección de datos para un subconjunto, o una muestra, de una población finita.
- Error de muestreo: diferencia entre la estimación de la encuesta y el verdadero parámetro de la población.
 - Cuanto mayor sea el error de muestreo, menor será la precisión.

Intervalo de confianza

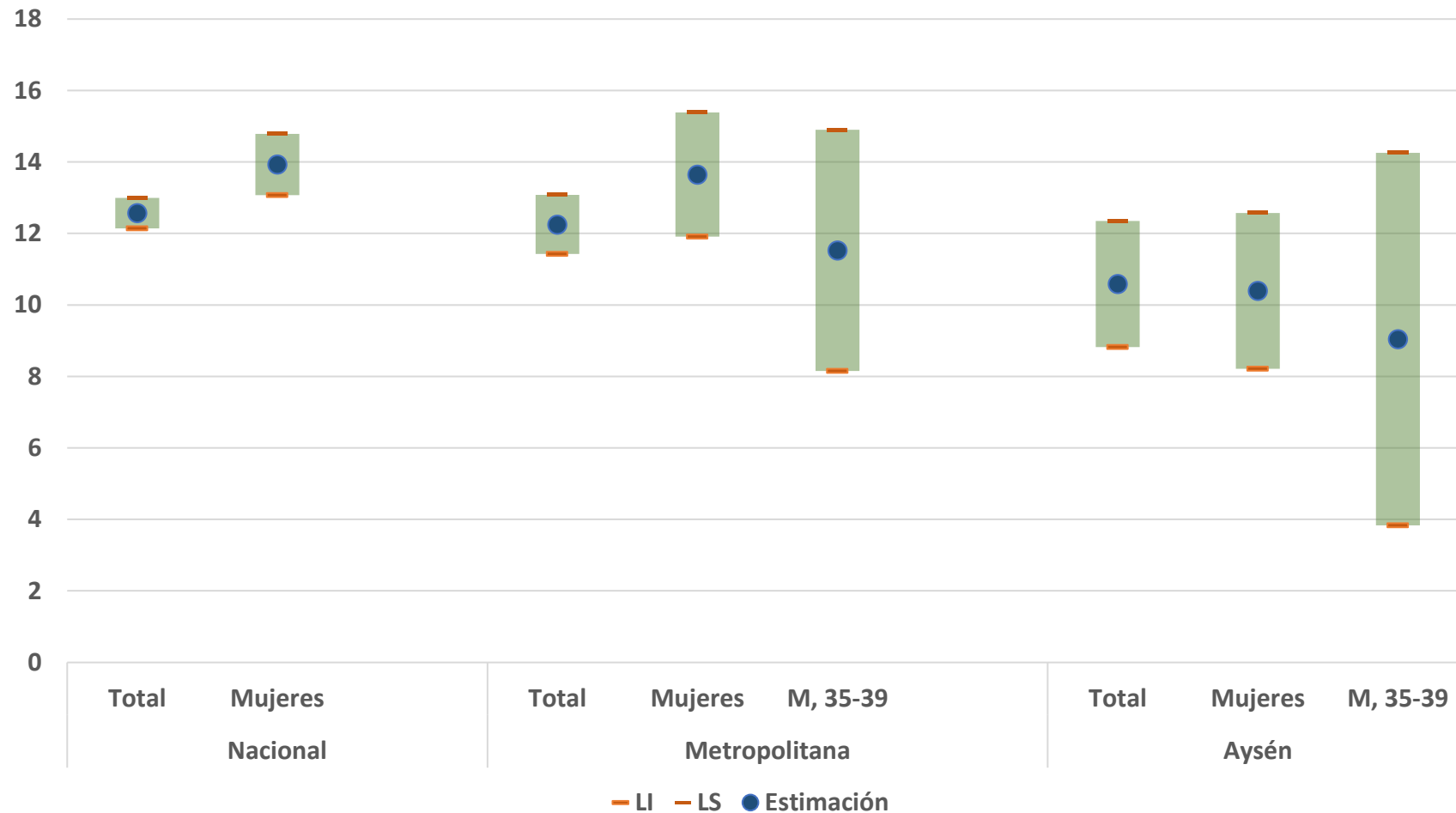
- El tamaño del error de muestreo se puede representar mediante el intervalo de confianza.
- Indica el posible rango de valores en el que es más probable que se encuentre el valor real del parámetro
 - “más probable” suele significar una probabilidad del 90%, 95% o 99%
- Un intervalo de confianza del 95% para un parámetro de interés (θ) viene dado por la siguiente expresión:

$$\left(\hat{\theta} - t_{0.975,gl} \times se(\hat{\theta}) \quad , \quad \hat{\theta} + t_{0.975,gl} \times se(\hat{\theta}) \right)$$

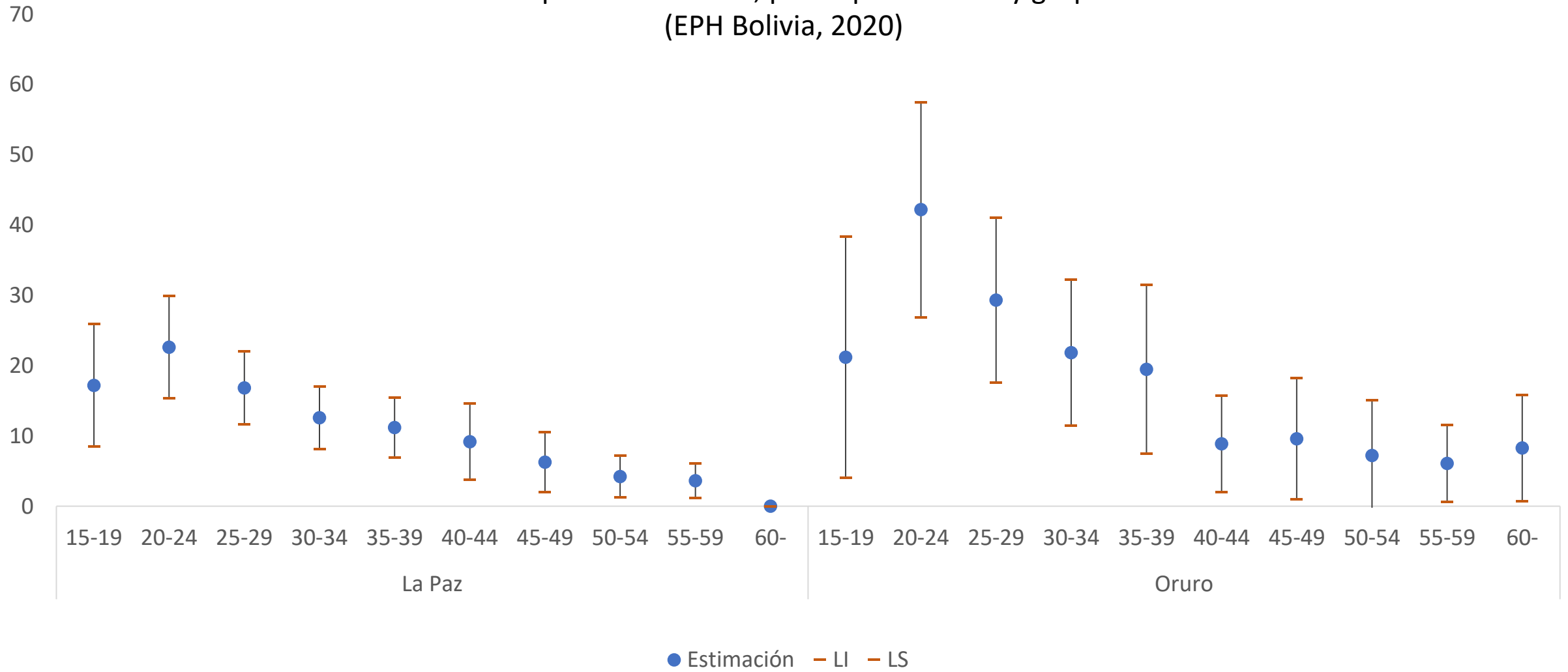
- Donde $\hat{\theta}$ es un estimador para el parámetro de interés, $t_{0.975,gl}$ es el percentil 0,975 de una distribución t-student con gl grados de libertad (UPM – estratos) y $se(\hat{\theta})$ es el error estándar.
- Los intervalos de confianza nos permiten inferir la precisión de un estimador

Algunos ejemplos

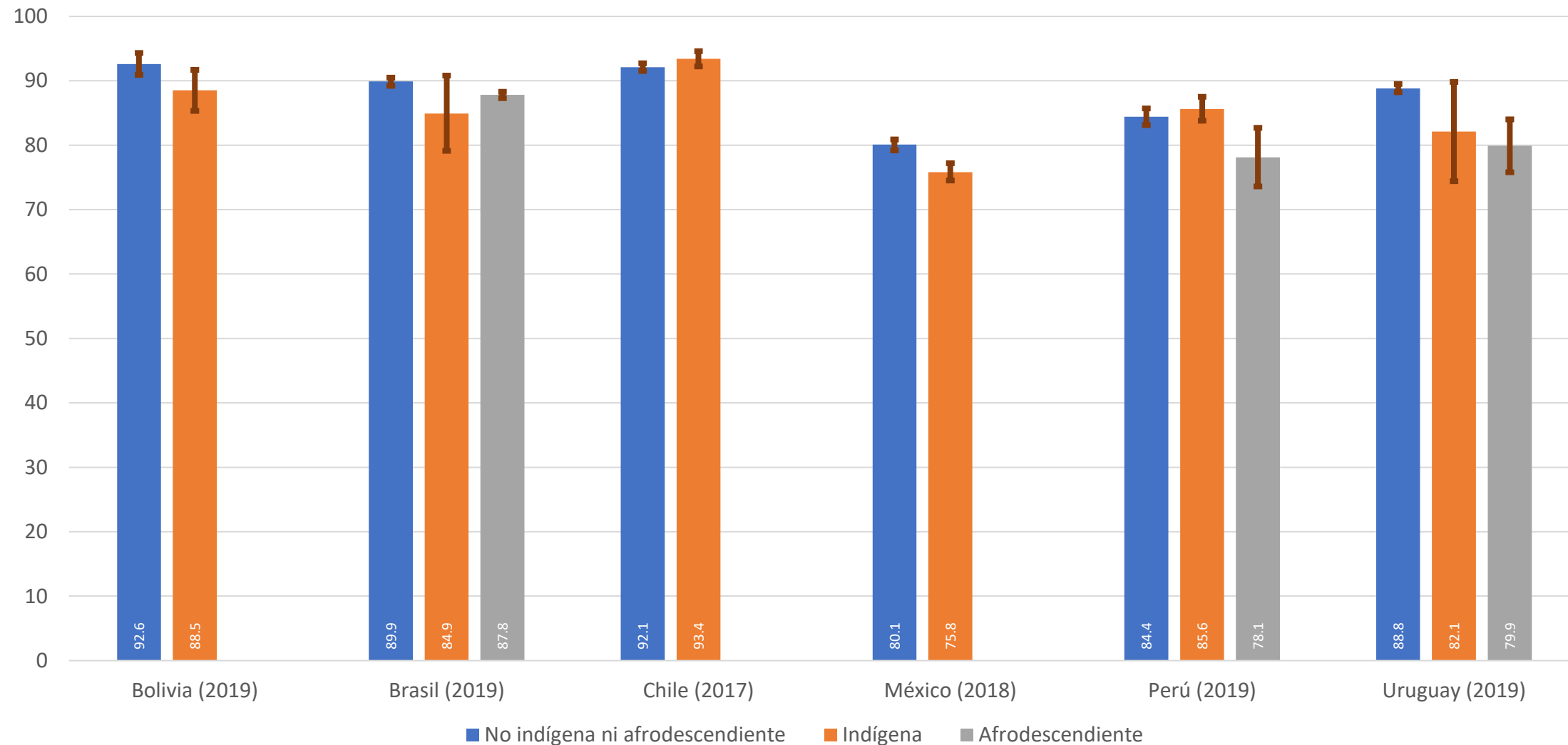
Tasa de desocupación, por regiones, sexo y grupo de edad (CASEN 2020)



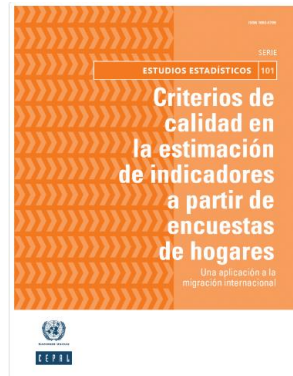
Tasa de desocupación femenina, por departamento y grupo de edad (EPH Bolivia, 2020)



Tasa de asistencia a la educación secundaria, población de 12 a 18 años de edad, por grupo étnico o raza



Criterios para evaluar la precisión de los estimadores



- Las Oficinas de Estadística aplican criterios heterogéneos para evaluar la precisión de los estimadores y determinar si deben ser publicados.
- Para avanzar hacia un enfoque más estandarizado:
 - CEPAL publicó en 2020 un [conjunto de criterios para evaluar la calidad de la estimación de indicadores mediante encuestas de hogares](#).
 - INE Chile implementó la [librería “calidad” en R](#) para evaluar la precisión de las estimaciones.
 - La CEA-CEPAL cuenta con un [Grupo de Trabajo](#) sobre “Recomendaciones para el análisis de la calidad de las encuestas de hogares” (2021-2023)
 - Sistematizar prácticas actuales
 - Procedimiento estandarizado para evaluar calidad y precisión
 - Procesos de anonimización y efecto sobre el error de muestreo



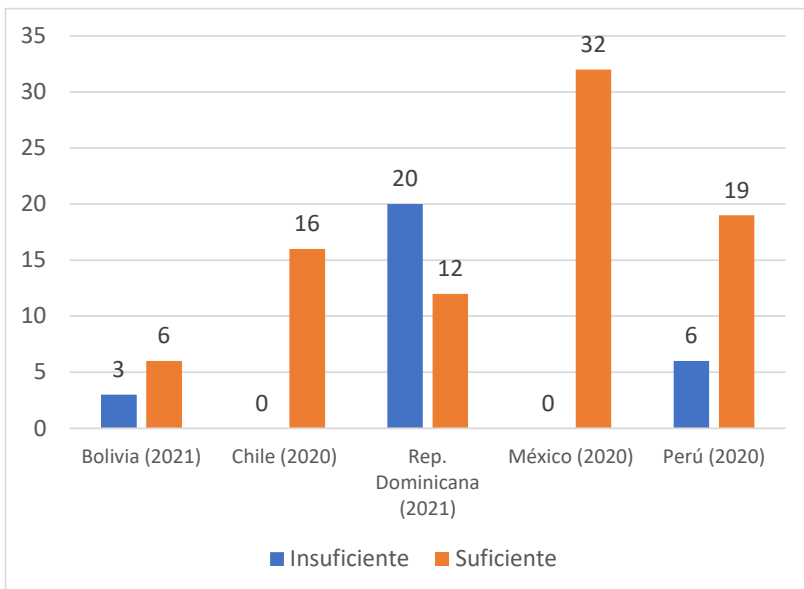
Criterios para evaluar la precisión de los estimadores

- Coeficiente de variación
 - Medida relativa del error de estimación
 - Generalmente usado por las ONE para evaluar la precisión (valores máximos entre 10% y 30%)
- Coeficiente de variación logarítmico
 - Transformación logarítmica sobre la proporción, para evitar coeficientes de variación artificialmente altos en estimaciones cercanas a cero.
- Tamaño de muestra
 - Número de casos en la muestra correspondientes a la subpoblación bajo estudio.
 - A mayor tamaño de muestra, menor error estándar.
- Tamaño de muestra efectivo
 - Tamaño de muestra disponible, corregido por el efecto de diseño.
- Conteo de casos no ponderado
 - Número de casos en la muestra afectados por el fenómeno de estudio.
- Grados de libertad
 - Número de UPM menos número de estratos

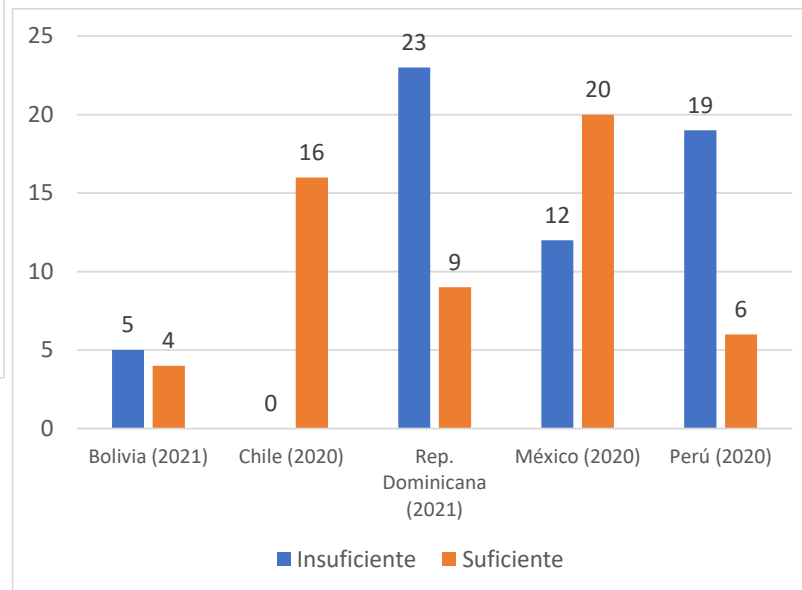
$$cve(\hat{\theta}) = \frac{se(\hat{\theta})}{\hat{\theta}}$$

$$n_{efectivo} = \frac{n}{Deff}$$

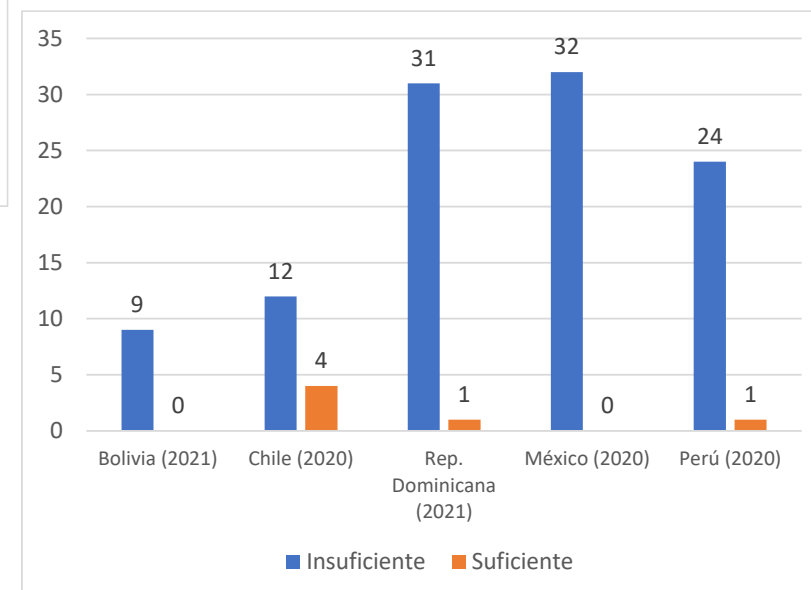
Ejemplo: Número de DAM según calidad de las estimaciones en 5 países



Tasa de desocupación, 15 y más años de edad



Tasa de desocupación femenina, 15 y más años de edad

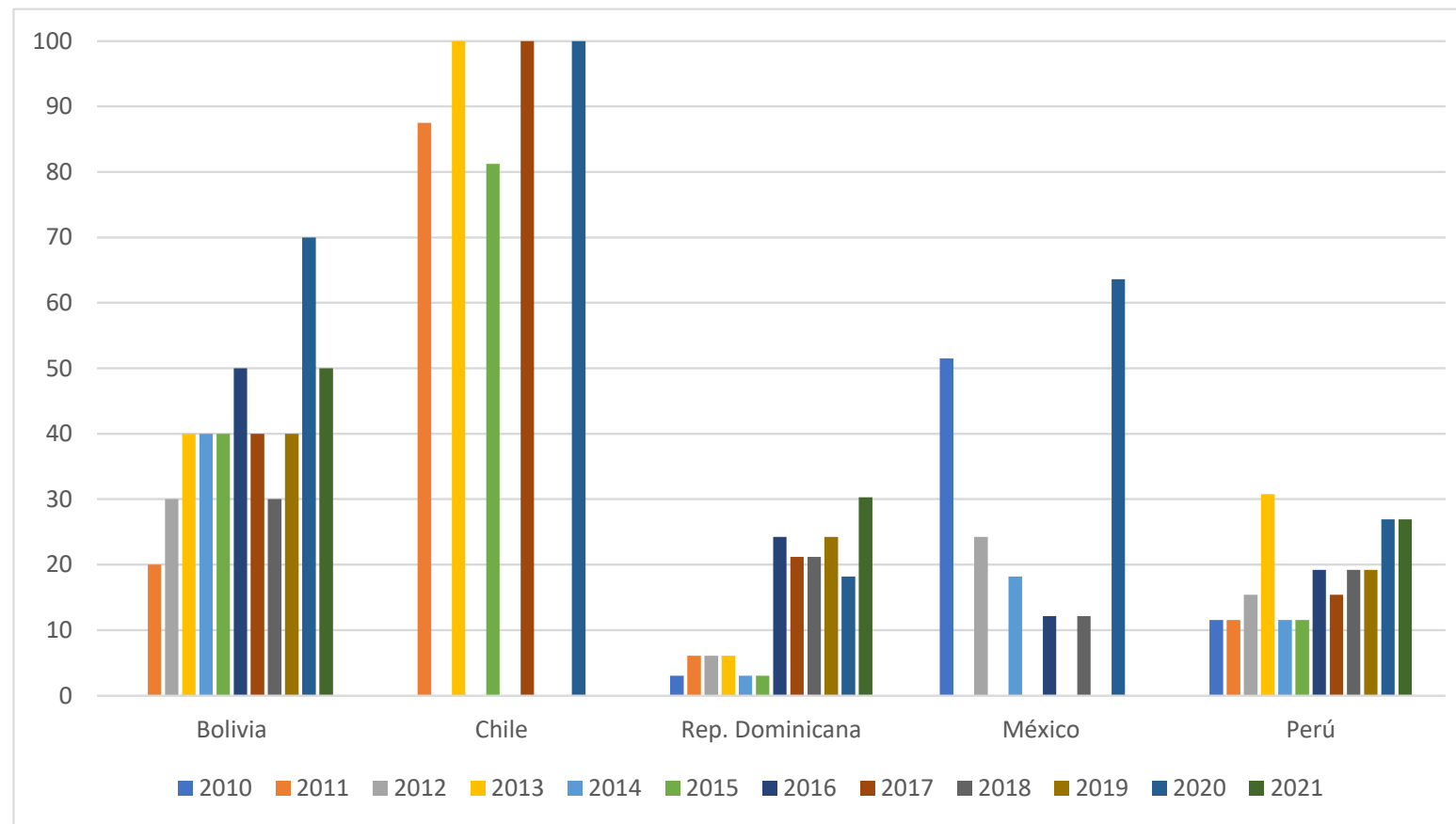


Tasa de desocupación femenina, 25 a 29 años de edad

Nota: DAM = División Administrativa Mayor

La disponibilidad de información con suficiente precisión puede variar a lo largo del tiempo, según los cambios en el diseño y características de las encuestas.

Porcentaje de DAM con datos válidos sobre la tasa de desocupación femenina



Estimaciones de áreas pequeñas (SAE) para la desagregación de datos

- Modelos estadísticos que permiten aprovechar la capacidad de desagregación de otras fuentes de datos:
 - Censos de población
 - Registros administrativos
 - Otras (datos satelitales, otras encuestas, etc.)
- Posibles ventajas:
 - Generar información a niveles que no es posible con los estimadores directos
 - Lograr mayores niveles de precisión que los estimadores directos
 - Generar estimaciones para dominios sin muestra en la encuesta
- “Área pequeña”:
 - Dominio para el cual no es posible obtener estimaciones directas confiables
 - Suelen ser dominios de estudio no planificados, con un tamaño de muestra esperado aleatorio
 - La subpoblación de interés puede ser un área geográfica o un subgrupo socioeconómico.

¿Cómo funciona SAE?

- 1. Identificar las variables auxiliares \mathbf{x} que están disponibles en la encuesta y en la fuente de datos complementaria.
- 2. Estimar un modelo para predecir la variable de interés \mathbf{y} , utilizando las variables auxiliares identificadas \mathbf{x} , con la encuesta.
- 3. Aplicar los parámetros estimados a la fuente complementaria, para predecir la variable de interés \mathbf{y} al nivel de desagregación deseado.
- 4. Medir el error de estimación y evaluar la confiabilidad de los resultados.

Modelos de área y modelos de unidad

- Modelos a nivel de área: predicen la variable de interés y a partir de variables auxiliares x agregadas a nivel de área (ej. comuna).
- Modelos a nivel de unidad: predicen la variable y usando variables auxiliares x para cada individuo y luego se agregan al nivel de área deseado.
- Los modelos a nivel de área también se pueden usar con datos a nivel de unidad, si la información se resume en el nivel de área apropiado.
- Los modelos a nivel de unidad pueden incorporar variables a nivel de área.

Comentarios finales

- La demanda por información desagregada suele exceder la capacidad de las encuestas para producir estimaciones precisas.
- Los métodos SAE ofrecen una forma práctica de producir desagregaciones que no son alcanzables mediante una estimación directa.
- Todos los métodos SAE requieren datos auxiliares al nivel del área pequeña de interés.
- Algunos aspectos a tener en consideración:
 - Obtener variables auxiliares con alto grado de asociación con la variable de interés.
 - Es importante estimar las medidas de error para SAE. Pero también tener en consideración que pueden existir otras fuentes de error.
 - Puede no ser posible usar las mismas fuentes de datos y modelo estadístico a lo largo del tiempo.

Gracias por su atención!

Xavier Mancero
División de Estadísticas
CEPAL

