

# Identificación de potencial sesgo por no respuesta en encuestas de hogares del nordeste argentino sobre el ingreso familiar y el nivel de pobreza

Maestría en Estadística Aplicada  
Universidad Nacional de Córdoba

Tesista	Lic. Celine Iliana Cabás
Directora	Dra. Patricia Caro
Co-Director	Dr. Carlos Matías Hisgen

Agosto 2024

## 1. Introducción

En el campo estadístico del muestreo, la no respuesta es una de las principales problemáticas a abordar debido a que si no se comporta de manera aleatoria puede introducir sesgos en las distribuciones de las variables que nos interesan medir y, por lo tanto, en las conclusiones que saquemos de testear hipótesis a partir de ellas.

”Si la decisión de responder depende estadísticamente de las variables bajo investigación, entonces la submuestra de encuestados no reflejará con precisión la distribución real de las variables de interés en la población y esto, a su vez, dará como resultado inferencias basadas en muestras sistemáticamente sesgadas”. [1]

Con respecto a esta problemática, un producto estadístico muy estudiado son las encuestas de hogares que relevan los países para medir diversos indicadores económicos y sociales. Entre sus finalidades, una de ellas es poder conocer la distribución del ingreso de las familias e individuos de cierto entorno geográfico y poder realizar inferencia sobre ellas.

A su vez, dentro de las medidas más relevantes calculadas a partir del ingreso, se encuentra el nivel de pobreza monetaria para las distintas unidades geográficas que conforman la muestra. Este indicador constituye una de las medidas centrales a tener en cuenta cuando se trata de implementación de políticas públicas.

El presente trabajo tiene como objetivo identificar si existe presencia de sesgo por no respuesta unitaria en el ingreso familiar medido por encuestas de hogares relevadas en nordeste argentino. El estudio se centra en los aglomerados urbanos de Gran Resistencia, Corrientes, Formosa y Posadas para los años 2018-2022 con periodicidad trimestral.

De los distintos enfoques que existen para abordar el problema de la no respuesta, este proyecto pretende estudiar el caso particular de la no respuesta 'unitaria' que hace referencia a la pérdida de observaciones por entrevistas no realizadas. Se trabajará con modelos que determinen la potencial presencia de sesgo por no respuesta y, en caso de existir, se propondrán métodos de reponderación de los datos para minimizar el sesgo.

En primera instancia, se pretende identificar si la cantidad total de entrevistas realizadas por un mismo hogar a lo largo del período depende estadísticamente o no del ingreso per cápita familiar y de otras variables que reflejen su condición socioeconómica. Se plantea un modelo logístico multinomial con variable de respuesta ordinal para estudiar si la predisposición a responder depende significativamente de las variables bajo estudio. [2] [1]

En el caso de probar la presencia potencial de sesgo, en segunda instancia se pretende trabajar con métodos de reponderación de los datos que corrijan por no respuesta. Principalmente, aquellos que utilicen la probabilidad predicha de respuesta de los hogares.

En este trabajo, se plantea esta problemática en el marco de la Encuesta Permanente de Hogares (EPH) para los aglomerados urbanos del nordeste argentino (NEA). Específicamente en lo que refiere al ingreso per cápita familiar y, por ende, al nivel de pobreza de los hogares. La selección de la región NEA viene justificada por el interés de estudiar si las marcadas diferencias entre el nivel de pobreza monetaria de Gran Resistencia y los demás aglomerados urbanos de la región puede verse justificada por diferencias en sus estructuras de no respuesta.

## 2. Antecedentes

En la literatura se han propuesto y estudiado procedimientos para minimizar el sesgo por no respuesta en las distintas etapas de una encuesta. Ya sea respecto a cómo tratar la no respuesta antes y durante la recolección de datos, previendo cuestiones vinculadas al contacto con la persona encuestada, como también métodos que evalúan sesgos potenciales posterior a la recolección y maneras de corregirlos [3] [4]. En esta investigación se trabajará en el segundo aspecto, ya que se emplearán datos secundarios relevados por programas nacionales.

El problema de la no respuesta puede entenderse principalmente de dos maneras. Como “no respuesta al ítem”, cuando el individuo no responde a un ítem en particular pero sí participa de la encuesta, o como “no respuesta unitaria”, donde el individuo directamente no realiza la entrevista [1].

En estudios basados en muestras, puede darse que la no respuesta no se comporte de manera aleatoria. Es decir, cuando las personas que sí respondieron a la encuesta difieren significativamente de aquellas que no lo hicieron, debe tenerse cuidado respecto a la validez de la inferencia porque se estaría trabajando con distribuciones sesgadas [5].

Entonces, ¿cómo detectar el sesgo por no respuesta? Existen diversas alternativas como el test de Little para determinar si el patrón de la no respuesta sigue un comportamiento completamente aleatorio [6]. Otra alternativa son los test chi-cuadrado de contraste de independencia, o bien modelos de regresión logísticos de respuesta binaria o politómica que permiten testear si existe relación entre ciertas variables de interés y la propensión a responder de los hogares.

Una vez detectado que la no respuesta no sigue un comportamiento aleatorio en la muestra, la siguiente cuestión a resolver es cómo lidiar con la no respuesta ¿de qué manera corregir el sesgo? Existen técnicas de imputación para sustituir individuos no encuestados, por ejemplo dividir la muestra en subclases e imputar un individuo similar a la misma subclase. Además, se desarrollaron métodos basados en variables auxiliares que son ampliamente utilizados pero tienen la limitación de que requieren información de la población. Estas variables deben estar disponible para los respondentes y al mismo tiempo conocer algún tipo de información sobre la distribución de esta variable en la población. [2]

Otra alternativa es el método para corregir la no respuesta mediante el ajuste del factor de expansión basada en la probabilidad de respuesta predicha, ampliamente conocido como Propensity Score Adjustment (PSA). Estos métodos parten desde modelos logísticos de respuesta binaria o modelos multinomiales de respuesta politómica hasta modelos más modernos vinculados a la literatura de Machine Learning (DecisionTree, Random Forest, XGBoost) que tratan de predecir la no respuesta para incluirla en la corrección del ponderador. [7] [2]

Una vez calculada la probabilidad, el método más usual es el de la ponderación inversa de la probabilidad que recalibra los pesos de los hogares con baja probabilidad de respuesta mejorando su representatividad en la muestra. [7]

### 3. Formulación del problema y objetivos

La **pregunta de investigación** planteada en este trabajo es la siguiente:

¿Qué efectos tiene la presencia de sesgo por no respuesta en encuestas de hogares del nordeste argentino sobre la estimación del ingreso familiar y el nivel de pobreza?

El **objetivo general** es identificar la potencial presencia de sesgo por no respuesta en encuestas de hogares del nordeste argentino y sus efectos sobre el ingreso familiar y el nivel de pobreza durante el período 2018-2022.

Los **objetivos específicos** son:

- Comparar las estructuras de respuesta de la encuesta de hogares entre los aglomerados urbanos del nordeste argentino para el período 2018-2022.
- Comprobar si la predisposición a responder por parte de los hogares depende significativamente del ingreso familiar u otras variables en los distintos aglomerados urbanos del nordeste argentino.
- Comparar modelos para predecir la probabilidad por parte de los hogares de responder de manera completa el esquema de entrevistas de la encuesta.
- Plantear una corrección del sesgo por no respuesta basada en la probabilidad predicha de responder mediante técnicas de ponderación de los datos.
- Contrastar las distribuciones de ingresos y el nivel de pobreza estimado antes y después de la corrección por no respuesta.

### 4. Fuentes de información

La principal fuente de información de este trabajo es la Encuesta Permanente de Hogares (EPH) relevada por el Instituto Nacional de Estadísticas y Censos (INDEC) de la República Argentina, que se constituye como un programa nacional de relevamiento de información sobre trabajo e ingresos de los hogares y es empleada para el cálculo de la pobreza. Se trabajará con las bases de datos por hogar e individuo de periodicidad trimestral para el período que comprende el primer trimestre 2018 al cuarto trimestre 2022.

Se complementa su uso con el Índice de Precios al Consumidor (IPC) para deflactar los ingresos y poder estudiar los efectos constantes durante el período. Además, se trabajó también con la valorización de la canasta básica total (CBT) para la determinación de la condición de pobreza o no del hogar por período. Ambas variables definidas para la región nordeste de Argentina.

La muestra sigue un esquema de rotación trimestral de las viviendas por área muestral. Es decir que, idealmente, una misma vivienda debe ser encuestada dos trimestres consecutivos, descansar los dos trimestres subsiguientes y volver a ser encuestada dos trimestres consecutivos más para garantizar la estructura de panel de datos que caracteriza a la encuesta. Sin embargo, debido a la no respuesta no todos los hogares completan el esquema lo cual nos sirve como variable proxy para medir la tendencia a responder que tienen las viviendas que son encuestadas.

Dado que los datos públicos de la encuesta incluyen únicamente los casos con entrevista realizada, para poder estudiar la probabilidad de respuesta de los hogares trabajaremos con una variable que cuenta el número de entrevistas realizadas en todo el período considerado. Es decir, los hogares que sólo contestaron un trimestre de los cuatro que corresponden tienen baja propensión a responder mientras que los que contestaron las cuatro veces tienen alta propensión a responder. De esta manera, podremos evaluar si existen diferencias sistemáticas en la probabilidad de responder más o menos veces según el nivel de ingreso y otras variables socioeconómicas.

Cuadro 1: Descripción de variables a utilizar en EPH individual y hogar.

Variable	Descripción
<b>Identificación</b>	
CODUSU	Código de identificación de la vivienda
NRO_HOGAR	Código de identificación del hogar
REGION	Código de región geográfica
AGLOMERADO	Código de aglomerado urbano
ANO4	Año de relevamiento
TRIMESTRE	Trimestre de relevamiento
<b>Base individual</b>	
CH03	Relación de parentesco (Jefe de hogar=1)
CH04	Sexo
CH06	Edad
NIVEL_ED	Nivel educativo
ESTADO	Condición de actividad
CAT_OCUP	Categoría ocupacional
<b>Base hogar</b>	
IV1	Tipo de vivienda
IX_TOT	Cantidad de miembros del hogar
ITF	Ingreso total familiar
IPCF	Ingreso per cápita familiar
<b>Variable de respuesta</b>	
NRO_REP	Número de entrevistas realizadas en el período 2018-2022
<b>Ponderador</b>	
PONDIH	Ponderador del ITF y del IPCF

## 5. Metodología

En una primera etapa, se realizará un análisis descriptivo de la estructura de la muestra para los cuatro aglomerados urbanos que representan al nordeste argentino en la encuesta (Gran Resistencia, Corrientes, Formosa y Posadas). Esto permitirá observar, a priori, si existen diferencias importantes en el número de entrevistas realizadas por hogar que pueda causar sesgos por no respuesta.

En una segunda etapa, se trabajará un modelo que explique el ingreso per cápita a precios constantes de los hogares en base a distintas características socioeconómicas con el objetivo de testear si, controlando por estas variables, el número de repeticiones tiene una relación significativa con el ingreso. Dada la estructura de panel de los datos se implementará un modelo mixto con ordenada aleatoria para controlar la heterogeneidad de las viviendas.

Siguiendo a [1], en una tercera etapa se ajustará un modelo lineal generalizado multinomial de respuesta politómica que intente explicar el número de repeticiones por hogar para testear la aleatoriedad teórica que debería tener esta variable respecto a características socioeconómicas de interés. En caso de identificar la presencia de sesgo porque la decisión de responder depende de las variables de interés, pasamos a la siguiente etapa.

Con el objetivo de predecir la probabilidad de respuesta de los hogares basado en la propensión que tengan a completar el esquema de la encuesta, se compararán modelos alternativos de predicción como ser el logístico, DecisionTree, Random Forest y XGBoost seleccionando el de mejor desempeño para esta aplicación en particular. Las probabilidades predichas por estos modelos se implementarán en la corrección de los factores de expansión de la encuesta.

Por último, para analizar el efecto de la no respuesta sobre el ingreso per cápita familiar y el nivel de pobreza por aglomerado urbano, se compararán las distribuciones de ingresos pre y post corrección. Lo cual implicará el análisis principalmente de las colas de la distribución que representan los percentiles de más bajos y mayores ingresos.

## 6. Resultados esperados

### Referencias

- [1] Martin Ravallion Anton Korinek Johan A. Mistiaen. “An econometric method of correcting for unit nonresponse bias in surveys”. En: *Journal of Econometrics* 136 (2007), págs. 213-235.
- [2] Fannie Cobben Jelke Bethlehem y Barry Schouten. *Handbook of nonresponse in household surveys*. John Wiley Sons, Inc., Hoboken, New Jersey, 2011.
- [3] Wendy Van de Kerckhove Thomas Krenzke y Leyla Mohadjer Westat. “Identifying and Reducing Nonresponse Bias throughout the Survey Process”. En: *ASA Section on Survey Research Methods* (2005).
- [4] Thomas Krenzke Wendy Van de Kerckhove y Leyla Mohadjer. “Approaches to a Nonresponse Bias Analysis in an Adult Literacy Survey”. En: *ASA Section on Survey Research Methods* ().
- [5] Chiao-chih Chang y Ferry Butar Butar. “Weighting Methods in Survey Sampling”. En: *Section on Survey Research Methods* (2012).
- [6] Leonardo González Allendes. “Propuesta de tratamiento de la no respuesta parcial para la medición de la Pobreza Multidimensional en Chile”. Tesis de mtría. Universidad de Chile, 2019.
- [7] María del Mar Rueda Ramón Ferri-García Jorge L. Rueda-Sánchez y Beatriz Cobo. “Estimating response propensities in nonprobability surveys using machine learning weighted models”. En: *Mathematics and Computers in Simulation* (2024).