



UNIVERSITAT OBERTA DE CATALUNYA (UOC)  
MÁSTER UNIVERSITARIO EN CIENCIA DE DATOS (*Data Science*)

## TRABAJO FINAL DE MÁSTER

ÁREA: INFORMÁTICA, MULTIMEDIA Y TELECOMUNICACIÓN

# Predicción de panel attrition con Machine Learning: El caso de la Encuesta Financiera de las Familias

---

Autor: Carlos Luis Gento de Celis

Tutor: Jordi Escayola Mansilla

Profesor: Antonio Lozano Bagén

---

Madrid, 22 de diciembre de 2023



# Créditos/Copyright



Esta obra está sujeta a una licencia de Reconocimiento - NoComercial - SinObraDerivada  
3.0 España de Creative Commons.



# FICHA DEL TRABAJO FINAL

Título del trabajo:	Predecir panel attrition con Machine Learning: Un análisis con la Encuesta Financiera de las Familias
Nombre del autor:	Carlos Luis Gento de Celis
Nombre del colaborador/a docente:	Jordi Escayola Mansilla
Nombre del PRA:	Antonio Lozano Bagén
Fecha de entrega (mm/aaaa):	10/2023
Titulación o programa:	Máster Universitario en Ciencia de Datos
Área del Trabajo Final:	Informática, Multimedia y Telecomunicación
Idioma del trabajo:	Español
Palabras clave	predictive models, machine learning, panel attrition



# Resumen

La Encuesta Financiera de las Familias (EFF) es una encuesta bienal cuyo objetivo es recoger información sobre la situación económico-financiera de los hogares que residen en España, y su evolución a lo largo del tiempo. Para ello, los hogares seleccionados pueden participar en hasta cuatro olas consecutivas de la encuesta. Sin embargo, hay hogares que abandonan el estudio antes de tiempo. Aunque estos abandonos no interrumpen el estudio, sí podrían afectar a sus resultados si el número de abandonos es demasiado grande, o si se concentra en colectivos específicos de la población. Es importante analizar las causas de estos abandonos y desarrollar herramientas que puedan evitarlos.

Este documento, en primer lugar, analiza las características de los hogares que participaron en la Encuesta Financiera de las Familias (EFF) en sus ediciones de 2017 y 2020, y si participaron en las olas de 2020 y de 2022. A continuación, plantea una serie de modelos de predicción basados en métodos de Machine Learning y evalúa su capacidad para predecir si un hogar abandonará la EFF en 2020 o en 2022. Finalmente, interpreta los resultados del modelo que mejor ha funcionado.

**Palabras clave:** predictive models, panel attrition, machine learning, household surveys





# Abstract

The Spanish Survey of Household Finances (EFF) is a biennial survey whose goal is to collect information about the economical and financial situation of households in Spain, and its evolution through time. To do so, selected households may take part in up to four consecutive waves of this survey. However, some households cease their participation prematurely. Although withdrawals do not interrupt the research study, they might affect its results if the number of withdrawals is too high, or if it is more likely to happen for certain groups of people. It is important to analyse the causes of this phenomenon and develop tools to prevent it.

Firstly, this paper analyses the characteristics of households that responded to the EFF in its editions of 2017 and 2020, and their participation during the 2020 or 2022 waves. Next, a series of predictive models based on machine learning methods are considered and evaluated for the exercise of predicting panel attrition in the EFF. Finally, it interprets the results of the most successful model.

**Key words:** predictive models, panel attrition, machine learning, longitudinal surveys, household surveys



# Índice general

Resumen	v
Abstract	vii
Índice	ix
<b>1. Introducción</b>	<b>3</b>
<b>2. Objetivos del proyecto, planificación y motivación personal</b>	<b>7</b>
2.1. Objetivos del proyecto . . . . .	7
2.2. Planificación del proyecto . . . . .	8
2.3. Motivación personal . . . . .	8
<b>3. Panel attrition: causas, presencia en la EFF y predicción</b>	<b>11</b>
3.1. Causas del panel attrition y cómo reducir su impacto durante la recolección de datos . . . . .	11
3.2. Falta de respuesta y panel attrition en la Encuesta Financiera de las Familias . .	12
3.3. Predicción de Panel Attrition . . . . .	14
<b>4. Datos y metodología</b>	<b>17</b>
4.1. Datos: La Encuesta Financiera de las Familias . . . . .	17
4.2. Metodología . . . . .	20
4.2.1. Recopilación de datos . . . . .	21
4.2.2. Algoritmos de Machine Learning: Entrenamiento, validación y test . . . .	21
<b>5. Resultados y posibles próximos pasos</b>	<b>25</b>
5.1. Análisis exploratorio . . . . .	26
5.2. Evaluación de modelos . . . . .	28
<b>6. Conclusiones y futuras líneas de trabajo</b>	<b>33</b>





# Capítulo 1

## Introducción

Las encuestas son una forma de recolección de datos que se fundamenta en plantear preguntas relevantes a una muestra de unidades muestrales, y utilizar sus respuestas para entender a una población en su conjunto. Desde una perspectiva temporal, hay dos tipos de encuestas: las de sección cruzada, y las longitudinales o panel. Las encuestas de sección cruzada realizan sus preguntas a las unidades muestrales en un momento concreto del tiempo, mientras que las longitudinales plantean sus preguntas de manera repetida a las mismas unidades muestrales durante un período de tiempo (semanas, meses, años...). En este documento nos vamos a centrar en las encuestas longitudinales.

La dimensión temporal de las encuestas longitudinales las convierten en una herramienta muy útil para poder analizar relaciones causales, ya que recogen cambios de opiniones, comportamientos o estados de los mismos encuestados panel a lo largo del tiempo. Sin embargo, la calidad de esos análisis depende de la cooperación exitosa y continuada de dichos encuestados durante las sucesivas ediciones u olas de la encuesta. El abandono prematuro y acumulado en el tiempo de participantes en un panel se conoce como **Panel Attrition** ([Watson and Wooden \(2009\)](#)).

Conceptualmente, las causas del Panel Attrition pueden clasificarse en tres categorías secuenciales: no-localización, no-contacto y no-cooperación ([Lepkowski et al. \(2002\)](#)). La no-localización se refiere a no conseguir localizar a los panelistas durante el período de campo. Esto suele ocurrir por cambios o errores en la información de contacto recogida durante la ola anterior (direcciones de residencia, teléfonos o email). Si un participante panel cambia de residencia o cambia de número de teléfono, es más costoso o incluso inviable localizarle.

La segunda causa es el no-contacto, que es el hecho de, después de localizar exitosamente al panelista, no conseguir establecer un contacto. Para que esto suceda, es necesario que el intento de contacto por parte de los encuestadores coincida temporalmente con la disponibilidad del panelista. Esto depende completamente del método de recolección de los datos. Por ejemplo,

en una entrevista personal, el panelista debe estar en su residencia justo en el mismo momento en el que el entrevistador hace la visita al hogar.

Finalmente, después de establecer el contacto, es necesario convencer al panelista para que vuelva a participar otra vez en la encuesta. La falta de cooperación en encuestas en general se ha analizado bastante en la literatura, y suelen destacar factores como las características demográficas de los encuestados o la temática de la encuesta (para un desarrollo más detallado, puede consultarse [Groves et al. \(1992\)](#)). Sin embargo, en el caso de las encuestas longitudinales, los encuestados tienen experiencia previa por haber participado anteriormente, y este hecho puede potenciar el efecto de factores como la duración de la entrevista, la carga cognitiva por pensar las respuestas o la fatiga por haber participado en varias ediciones anteriores ([Laurie et al. \(1999\)](#), [Watson and Wooden \(2009\)](#), [Lynn \(2018\)](#)).

Con respecto a las consecuencias del Panel Attrition, [Lynn \(2018\)](#) destaca dos principales problemas. Por un lado, si la tasa de attrition es alta, el tamaño de la muestra se reducirá drásticamente con el paso de las olas, lo que provocará que la precisión de los estimadores de la encuesta sea muy baja, y además limitará o incluso imposibilitará el análisis de subgrupos dentro de la muestra. Por otro lado, si el Panel Attrition no es aleatorio, y por tanto los panelistas que abandonan la encuesta son sistemáticamente diferentes a los que se mantienen, existe el riesgo de introducir un sesgo de no-respuesta en los estimadores.

Tradicionalmente, los efectos del Panel Attrition se han mitigado con el uso de métodos de imputación múltiple ([Rubin \(1987\)](#)), reponderando de pesos muestrales ([Groves et al. \(2009\)](#)) e introduciendo muestras de refresco para sustituir a las unidades muestrales perdidas ([Hirano et al. \(1998\)](#)). Pero en las últimas décadas se ha extendido el uso de los diseños adaptativos y reactivos (adaptive and responsive designs, [Groves and Heeringa \(2006\)](#), [Wagner \(2008\)](#), [Schouten et al. \(2017\)](#), [Tourangeau et al. \(2017\)](#)). La idea detrás de estos diseños se fundamenta en utilizar toda la información que se genera durante la elaboración de encuestas (respuestas al cuestionario, paradata u observaciones de los entrevistadores) para diseñar implementaciones informadas cuyo objetivo sea mejorar la calidad de los datos, reducir los costes, o ambos. Por ejemplo, se han utilizado para revisar incentivos a participar para grupos concretos de encuestados [Mcgonagle et al. \(2022\)](#) o revisar el orden de las preguntas ([Early et al. \(2017\)](#)). En este sentido, las encuestas longitudinales ofrecen una gran oportunidad para estos diseños porque contienen mucha información tanto de los trabajos de campo que se estén desarrollando, como los que se realizaron en olas anteriores. Por ejemplo, en [Kreuter and Müller \(2015\)](#) utilizan los registros de llamadas a panelistas en ediciones anteriores para optimizar las estrategias de contacto en las ediciones siguientes.

Dentro de este contexto de los diseños adaptativos y reactivos, en las últimas décadas se ha desarrollado el uso algoritmos de Machine Learning en la metodología de encuestas ([Buskirk](#)

---

et al. (2018), Kern et al. (2019)). Y de manera particular, han presentado resultados prometedores a la hora de predecir resultados de participación en los trabajos de campo, especialmente los modelos basados en árboles de decisión (Kern et al. (2019), Kern et al. (2021), Liu (2020)). En el contexto de Panel Attrition, en Beste et al. (2023) utilizaron información de ediciones pasadas de una encuesta a hogares en alemania para entrenar un Random Forest e identificar hogares panelistas con una baja propensión a participar. Luego, utilizaron esa información para crear un diseño experimental en la siguiente ola de la encuesta, en el cual se asignaba un incentivo monetario adicional a la mitad de esos hogares. Sus resultados mostraron incrementos en las tasas de respuesta de los hogares tratados, y animaron a sus responsables a seguir utilizando este diseño adaptativo en futuras ediciones.

El objetivo de este Trabajo de Fin de Máster es adaptar la implementación de machine learning vista en Beste et al. (2023) para predecir la participación de los hogares panel en el contexto de la Encuesta Financiera de las Familias (EFF). En cada uno de estos artículos se utiliza información sobre hogares que participaron en olas pasadas de dos encuestas, y se utiliza esa información para predecir una variable binaria de participación o no participación en una ola posterior. El rendimiento de todos los modelos se compara con un modelo de referencia, que es una regresión logística. En el caso de este proyecto, utilizamos datos de las olas 4, 5, 6 y 7 de la EFF (que se corresponden con los años 2011, 2014, 2017 y 2020) para entrenar varios modelos de machine learning. A continuación, utilizamos dichos modelos para predecir la participación de los hogares panel en la ola 8 (que se corresponde con el año 2022), y comparamos su rendimiento con respecto a un modelo de regresión logística.

El resto de este documento se organiza de la siguiente manera. En el siguiente capítulo se exponen los objetivos, la planificación y las motivaciones personales de este proyecto. En el tercer capítulo se presenta la EFF, los datos que se han utilizado y la metodología aplicada. En el cuarto capítulo se presentan los resultados de la implementación descrita en el tercer capítulo, junto con los desafíos encontrados durante el proceso. Finalmente, el último capítulo contiene las conclusiones obtenidas de los resultados del proyecto y reflexiones sobre cuáles podrían ser los próximos pasos para continuar con este proyecto en el futuro.





## Capítulo 2

# Objetivos del proyecto, planificación y motivación personal

### 2.1. Objetivos del proyecto

Este proyecto se contextualiza dentro del marco de la metodología de encuestas. En concreto, se centra en las encuestas longitudinales a hogares, y dentro de ese campo pone el foco en la predicción de la participación de los hogares en olas posteriores a su primera colaboración. La encuesta elegida para el proyecto es la Encuesta Financiera de las Familias (EFF), una encuesta de referencia para investigaciones sobre finanzas de los hogares.

El **objetivo principal** de este proyecto es desarrollar un modelo de predicción basado en métodos de machine learning que ayude a predecir si un hogar que ha participado en al menos una ola de la EFF volverá a hacerlo en olas posteriores.

Para poder desarrollar ese modelo, es necesario completar una serie de **objetivos secundarios**. Estos objetivos se agruparán en 5 fases:

1. Hacer una revisión del estado del arte sobre el Panel Attrition y las metodologías de Machine Learning aplicadas a su predicción.
2. Recolección de datos. Preferentemente, obtener un conjunto de datos que contenga:
  - a) Las respuestas de los hogares al cuestionario de la EFF.
  - b) Paradata sobre el proceso de creación de la encuesta.
  - c) Características de los entrevistadores.
3. Análisis exploratorio de los datos. Identificar patrones dentro de los datos que puedan estar relacionados con el panel attrition.

4. Preprocesar los datos para entrenar modelos de machine learning.
5. Modelado y evaluación de modelos de predicción basados en métodos de machine learning.
6. Interpretación del mejor modelo y redacción de conclusiones.

## 2.2. Planificación del proyecto

La planificación de este proyecto se va a dividir en 5 fases. La asignación temporal a cada una de ellas se ha realizado de acuerdo a la estructura de contenidos del Plan Docente y al calendario de la asignatura.

1. Fase 1: Contiene la definición del proyecto y la planificación del TFM. Abarca las dos primeras semanas del proyecto, del 27 de septiembre al 10 de octubre.
2. Fase 2: Contiene las tareas de la revisión de la literatura y la caracterización del panel attrition de la EFF. Abarcará desde el 11 de octubre al 23 de octubre.
3. Fase 3: Contiene las tareas de recolección de datos, análisis exploratorio de dichos datos, preprocesamiento para entrenar los modelos de machine learning, modelado y evaluación de los modelos de predicción, y la interpretación de los modelos y la redacción de conclusiones. Esta fase es la más duradera y abarcará desde el 17 de octubre al 19 de diciembre.
4. Fase 4: En esta fase se procederá a la redacción de la memoria del TFM y la preparación de una presentación audiovisual sobre el proyecto. Abarcará desde el 20 de diciembre al 18 de enero.
5. Fase 5: En esta fase final se procederá a la defensa del TFM. Esa etapa abarcará desde el 22 de enero hasta el 4 de febrero.

En la figura [2.1](#) puede observarse la distribución temporal y el tiempo dedicado a cada tarea en un diagrama de Gantt.

## 2.3. Motivación personal

Trabajo para el Banco de España, y formo parte del equipo que elabora la Encuesta Financiera de las Familias (EFF) desde principios de 2015. He participado en la elaboración de sus cuatro últimas ediciones (EFF2014, EFF2017, EFF2020 y EFF2022) y con los años he adquirido cada vez más interés por la metodología de encuestas y el potencial que tienen para recoger

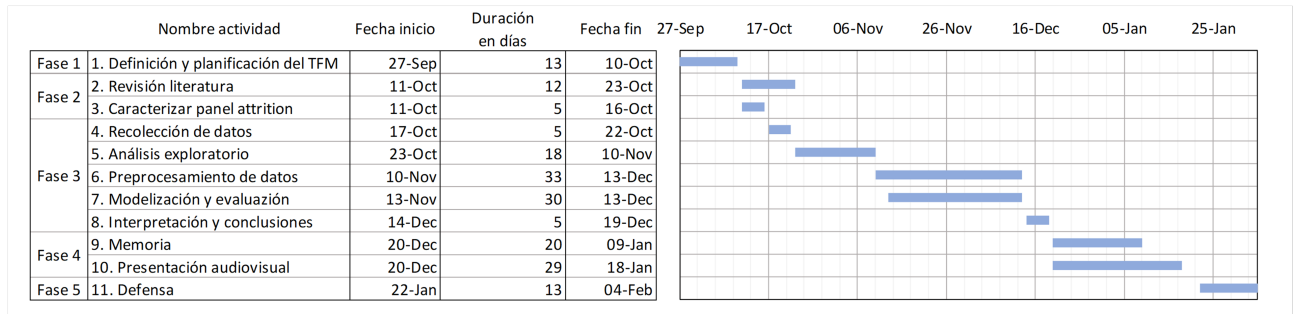


Figura 2.1: Planificación de las actividades del TFM

información sobre fenómenos que de otra manera serían difíciles de captar. La EFF es una fuente de información de referencia en el campo de las finanzas de los hogares, y el aumento del interés que se ha generado en los últimos años hace más importante realizar trabajos y esfuerzos para garantizar e incluso mejorar la calidad de sus datos. En ese sentido, la gran cantidad de parámetros que se generan durante cada ola ofrecen muchas oportunidades para aprender sobre todo el proceso, y también mejorarlo.

Por otro lado, en los últimos años ha aumentado la variedad de formatos en los que recibimos los datos, tomando una gran relevancia los audios de las entrevistas y los comentarios de texto escritos por los entrevistadores. Los métodos y las herramientas desarrolladas en el campo de la ciencia de datos abren un mundo de posibilidades para poder analizar toda esa información, y que hace apenas unos años no eramos capaces ni de imaginar. Es un buen momento y una gran oportunidad para trabajar en este campo.

Finalmente, sobre el Panel Attrition, considero que la etapa más importante de la elaboración de la EFF es el trabajo de campo. En ella se contacta a los hogares, se les convence para participar en la encuesta y se realizan las entrevistas. Marca el devenir las siguientes etapas, y también el nivel de calidad de los datos. Por esa razón, es importante conseguir la colaboración de los hogares. Especialmente la de los hogares panel, ya que representan una proporción importante de la muestra, y no convencerles puede llegar a ser muy costoso. Y es un área que en la EFF no se ha podido explorar hasta. Es una gran oportunidad para aprender.



## Capítulo 3

# Panel attrition: causas, presencia en la EFF y predicción

### 3.1. Causas del panel attrition y cómo reducir su impacto durante la recolección de datos

Conceptualmente, las causas del panel attrition pueden clasificarse en tres categorías condicionales ([Lepkowski et al. \(2002\)](#)): no-localización, no-contacto y no-cooperación.

La no-localización se refiere a no localizar exitosamente a un encuestado durante una ola posterior. Generalmente, esto se debe a cambios en la información de contacto (dirección de residencia, número teléfono, correo electrónico...) obtenida del participante durante la ola anterior ([Couper and Ofstedal \(2009\)](#)). Algunos factores que pueden contribuir al éxito o fracaso en la localización son el método de recolección de datos, la propensión a los cambios de localización de los encuestados entre diferentes olas, el tiempo transcurrido entre olas e incluso el presupuesto ([Lynn \(2009\)](#)). Por ejemplo, las encuestas con entrevistas cara a cara utilizan métodos de rastreo y búsqueda que suelen ofrecer altos índices de localización y cooperación ([De Leeuw et al. \(2005\)](#), [Couper and Ofstedal \(2009\)](#)), pero también requieren esfuerzos adicionales para localizar a los participantes panel que se hayan mudado, como por ejemplo reembolsar a los entrevistadores los gastos derivados del proceso de búsqueda.

Posteriormente, tras localizar al encuestado, es necesario establecer un contacto. En este paso el método de recolección de datos vuelve a ser muy importante. Por ejemplo, en una encuesta por correo postal o por email, el contacto depende de que el encuestado reciba dicho correo y además le preste atención. En cambio, en entrevistas presenciales en el hogar o entrevistas telefónicas, debe darse una coincidencia temporal entre la disponibilidad del entrevistado (está en casa o tiene su teléfono disponible) y el intento de contacto del entrevistador (realiza una visita

su casa o llama por teléfono). En este sentido, dos estrategias de contacto que han presentado buenos resultados han sido utilizar un número de intentos de contacto alto y diversificado en horarios (mañana, tarde, fines de semana...), y establecer períodos para realizar entrevistas lo suficientemente largos (Nicoletti and Peracchi (2005), Watson and Wooden (2009)). Además, las encuestas longitudinales ofrecen la ventaja de poder utilizar datos sobre el proceso de recolección en olas anteriores, y diseñar mejores estrategias de contacto o identificar casos en los que el contacto puede ser complicado (Calderwood et al. (2012), Lagorio et al. (2016)).

Finalmente, cuando se ha establecido contacto con el encuestado, éste puede cooperar de nuevo con los encuestadores, o por el contrario rechazar hacerlo. La falta de cooperación es una preocupación común en las encuestas, y hay bastante literatura sobre las motivaciones que puede haber detrás de los rechazos. Suelen destacar factores relacionados con las características de los participantes, como características sociodemográficas (edad, sexo, nivel de salud...) o el entorno (composición del hogar, del barrio de residencia...); con el diseño de la encuesta, como la temática de las preguntas o el método de recolección de datos; y con las características de los entrevistadores, como su apariencia o la experiencia previa realizando encuestas (Groves et al. (1992), O'Brien et al. (2006)). En el caso de las encuestas longitudinales, los encuestados ya tienen una experiencia previa sobre la encuesta que puede hacer potenciar el efecto de factores como la duración de la entrevista, la carga cognitiva que supone pensar las respuestas, la sensibilidad sobre la temática o la fatiga por haber participado ya en varias ediciones (Laurie et al. (1999), Watson and Wooden (2009), Lynn (2018)). Pero, al igual que con los contactos, las encuestas longitudinales también poseen información sobre el proceso de contacto, el desarrollo de la entrevista y las características de los encuestados en olas anteriores que pueden ayudar a diseñar incentivos que contribuyan a la cooperación (Laurie and Lynn (2009)) o descubrir la influencia de los entrevistadores para convencer a los encuestados, y lo importante que es su continuidad entre olas (Lynn et al. (2014)).

### 3.2. Falta de respuesta y panel attrition en la Encuesta Financiera de las Familias

La Encuesta Financiera de las Familias es una encuesta longitudinal en la que se entrevista a hogares residentes en España. Su primera edición se realizó en el año 2002, y se ha producido de manera trienal hasta el año 2020. Desde entonces, su producción es bienal, siendo la ola de 2020 la primera planificada para ser publicada con esta nueva frecuencia.

El objetivo de su primera edición era captar bien la distribución de la riqueza de los hogares españoles, por lo que el diseño de su muestra se fundamentó en un sobremuestreo de hogares

### 3.2. Falta de respuesta y panel attrition en la Encuesta Financiera de las Familias

con mayor nivel de riqueza<sup>1</sup> (Bover (2004)). Desde la siguiente edición (EFF2005) en adelante, además de mantener la representatividad de la muestra para el año de la ola correspondiente, se incluyó como objetivo adicional el incluir un componente longitudinal que consistía en volver a preguntar a todos los hogares que participaron en la ola anterior, y así poder hacer análisis de causalidad. Para conseguir ambos objetivos, en las ediciones de 2005, 2008 y 2011 se implementaron muestras de refresco por estrato de riqueza para complementar a la muestra panel que venía de olas anteriores (Bover (2008), Bover (2011), Bover et al. (2014)). Finalmente, en la edición de la EFF2014 se estableció el límite máximo de participación de los hogares en cuatro olas consecutivas, y desde entonces se eliminan de la muestra longitudinal a aquellos hogares que ya han participado en cuatro olas consecutivas (Bover et al. (2018), Barceló et al. (2021)).

El proceso de recopilación de los datos se basa en la realización de entrevistas personales, generalmente en la residencia de los hogares seleccionados. Los entrevistadores deben localizar la vivienda en la que reside el hogar, realizar varias visitas en diferentes combinaciones de horarios y días de la semana<sup>2</sup>, pedir su colaboración con la EFF y concertar una cita para realizar la entrevista<sup>3</sup>. Este procedimiento es el mismo tanto para la muestra de refresco como para la muestra panel. La participación es completamente voluntaria, la información es confidencial, y los hogares pueden decidir interrumpir la entrevista o pararla en cualquier momento.

El Banco de España incluye un análisis de la falta de respuesta y de panel attrition en los documentos metodológicos de cada ola de la EFF (Bover (2004), Bover (2008), Bover (2011), Bover et al. (2014), Bover et al. (2018), Barceló et al. (2021)). En estos análisis se muestran, diferenciando entre muestra panel y muestra no panel, el número de hogares contactados por tipo de respuesta (completas, rechazos, no contactados...), las tasas de cooperación según el estrato de riqueza de los hogares contactados, y los odd-ratios de dos modelos logit de cooperación frente a rechazo en el que se utilizan diferentes categorías recogidas por los entrevistadores sobre las condiciones de los edificios, del nivel socioeconómico del barrio, el tamaño del municipio de resi-

---

<sup>1</sup>El diseño del sobremuestreo de la EFF toma como referencia la [Survey of Consumer Finances \(SCF\)](#), que es la encuesta equivalente a la EFF en Estados Unidos, y que elabora la Reserva Federal. Lo ejecuta el Instituto Nacional de Estadística, con la colaboración de la Agencia Tributaria. El nivel de riqueza de los hogares se obtiene a partir de las declaraciones individuales más recientes en el impuesto sobre el patrimonio que hay en España (facilitado por la Agencia tributaria), y se definen unos estratos de riqueza a partir de los intervalos de la SCF y los percentiles de la distribución del nivel de riqueza de los hogares. Los hogares en los estratos más altos tienen mayor probabilidad de ser seleccionados. Para las regiones de País Vasco y Navarra no se realiza sobremuestreo porque no se dispone de información

<sup>2</sup>De manera general, unos días antes de la primera visita, la empresa de campo envía una carta al hogar, firmada por el gobernador del Banco de España, para informarle de que ha sido seleccionado para participar en la EFF, y que en los próximos días recibirá la visita de una persona para concertar una cita para una entrevista personal. A algunos hogares panel también se les contacta por teléfono si mostraron alguna preferencia por ser contactados de esa manera.

<sup>3</sup>Los hogares que aceptan participar y completan la entrevista reciben un obsequio por parte del Banco de España. Los hogares panel también reciben otro obsequio por haber participado en la edición anterior, independientemente de que accedan a participar de nuevo.



dencia en número de habitantes y la región de residencia. En general se observa que se contacta con más hogares no panel que con hogares panel, pero es un resultado esperable porque sus tasas de cooperación también son más bajas que las de la muestra panel. Esto se mantiene tanto para la totalidad de la muestra como diferenciando por estrato de riqueza. También se observa que para ambos tipos de muestra las tasas de cooperación tienden a ser menores a medida que aumenta el estrato de riqueza de los hogares. Con respecto a los modelos logit, los resultados difieren entre olas ya que en una ola son estimadores significativos, pero en otra no lo son. Aún así, algunas tendencias sí se han observado de manera repetida en la mayoría de olas son que la probabilidad de cooperar decrece cuando aumenta el tamaño del municipio<sup>4</sup>, que la probabilidad de no cooperar aumenta para barrios con menor nivel económico, y que existen diferencias en la cooperación entre regiones. Sin embargo, es posible que el efecto regional esté afectado por algún tipo de efecto de entrevistador, ya que los entrevistadores suelen hacer sus entrevistas en las mismas regiones ola tras ola.

Estos análisis que acabamos de describir muestran que existe no respuesta y panel attrition en la EFF. Pero se limitan a hacer una comparación de cada edición con la inmediatamente anterior. No se observa cuántos hogares suelen completar las cuatro olas para las que han sido seleccionadas, ni la distribución de los abandonos con el paso de las olas. Por otro lado, a la hora de analizar la no respuesta de muestra panel no se utiliza información recogida durante la ola anterior que podría ser relevante para analizar el panel attrition, como por ejemplo las características de los hogares, el número de intentos de contacto en la edición anterior, o la duración de la entrevista de la ola anterior.

La EFF ofrece muchas oportunidades para poder investigar sobre la falta de respuesta y el panel attrition en esta encuesta, ya que hay muchos análisis que no se han podido realizar. Conocer las causas que pueden provocar el panel attrition en la EFF ayudará a desarrollar herramientas que ayuden a combatirlo, y con ello mejorar las tasas de respuesta de los hogares y la calidad de los resultados de la encuesta.

### 3.3. Predicción de Panel Attrition

La regla tradicional utilizada para diseñar encuestas ha sido la estandarización de todos los procesos y protocolos. Todas unidades muestrales deben ser tratadas de la misma manera. Con la excepción de las introducciones de los entrevistadores (Groves et al. (1992)), esta regla se mantuvo durante bastante tiempo. Pero el aumento de las reticencias de la población a participar en encuestas y los recortes presupuestarios llevó a buscar cómo mejorar la eficiencia

---

<sup>4</sup>Este efecto siempre es negativo para la muestra no panel en todas las olas, pero para la muestra panel hay algunas olas para las que su efecto no es significativo.

de los procesos de elaboración de encuestas, y empezaron a considerarse diseños adaptativos enfocados a subgrupos específicos de la muestra ([Groves and Heeringa \(2006\)](#), [Lynn \(2014\)](#), [Lynn \(2017\)](#)). En ese sentido, la predicción se ha convertido en una opción muy interesante para poder identificar anticipadamente a individuos que potencialmente podrían abandonar precipitadamente una encuesta longitudinal, y en los que los métodos de machine learning tienen un peso importante ya que no necesitan un conocimiento previo de las relaciones que se quieren estudiar, y suelen adaptarse bien a contextos en los que la relación entre la variable dependiente y sus predictores suele ser compleja y no lineal ([Buskirk et al. \(2018\)](#), [Kern et al. \(2019\)](#), [Kern et al. \(2021\)](#), [Jankowsky and Schroeders \(2022\)](#)).

Los estudios de modelos predictivos realizan una comparación de rendimiento entre modelos tradicionales utilizados para analizar panel attrition, casi siempre una regresión logística, y modelos basados en métodos de Machine Learning, como diferentes tipos de árboles de decisión o máquinas de soporte de vectores. [Kern et al. \(2019\)](#) y [Kern et al. \(2021\)](#) muestran casos en los que los modelos predictivos basados en machine learning, especialmente árboles de decisión, presentan resultados prometedores. Sin embargo, en [Jankowsky and Schroeders \(2022\)](#) apenas observan diferencias significativas entre los resultados de un modelo de regresión logística y los de un modelo GBM (Gradient Boosting Machine). La conclusión de ese artículo es que no necesariamente un modelo más complejo (y más complejo de entender) puede ser más adecuado para predecir panel attrition, y por tanto utilizarse para diseñar políticas adaptativas para reducirlo.

La intención de este proyecto es realizar un estudio similar a los propuestos por [Kern et al. \(2021\)](#) y [Jankowsky and Schroeders \(2022\)](#), y comprobar si modelos basados en métodos de machine learning pueden predecir adecuadamente el panel attrition en la EFF, y dar la posibilidad de contribuir al desarrollo de herramientas que ayuden a reducirlo.



# Capítulo 4

## Datos y metodología

### 4.1. Datos: La Encuesta Financiera de las Familias

Los datos que se van a utilizar en este proyecto provienen de la Encuesta Financiera de las Familias (EFF). La EFF es una encuesta oficial a hogares elaborada por el Banco de España y está incluida en el Plan Estadístico Nacional. Su primera edición se realizó en el año 2002, y se ha producido de manera trienal hasta el año 2020. Desde entonces, se produce de manera bienal.

El objetivo de la EFF es recabar información sobre las condiciones financieras de los hogares residentes en España. Es la única fuente estadística que permite relacionar información sobre activos, deudas, ingresos y gastos de los hogares españoles. Esto permite analizar las decisiones de inversión y financiación de las familias y conocer su situación patrimonial, y gracias a ello tener un mayor conocimiento de la economía española y ser un apoyo importante para el diseño de políticas públicas. Por nombrar algunos ejemplos de estudios realizados con la EFF, se ha utilizado para cuantificar el ahorro adicional generado para los partícipes en planes de pensiones de empresa ([Gómez and Villanueva \(2022\)](#)), para caracterizar cómo afectó la pandemia del Covid-19 a la situación patrimonial de los trabajadores más afectados por dicha crisis ([Alvar-gonzález et al. \(2020\)](#)) o para analizar las diferencias en aceptación y uso de tarjetas de crédito y banca online entre diferentes grupos de hogares desde el año 2002 ([Crespo et al. \(2023\)](#)).

El diseño de la muestra de la EFF tiene dos características importantes: un sobremuestreo de hogares ricos y un componente longitudinal o panel. El sobremuestreo de ricos<sup>1</sup> garantiza poder analizar con suficiente precisión el comportamiento de los hogares de la parte alta de la distribución de riqueza. Este detalle es importante porque la distribución de la riqueza entre

---

<sup>1</sup>La muestra de la EFF es seleccionada por el Instituto Nacional de Estadística, en colaboración con la Agencia Tributaria, a partir de las declaraciones individuales más recientes en el Impuesto sobre el Patrimonio. Los detalles sobre el muestreo de la EFF pueden consultarse en los documentos metodológicos disponibles en su [sitio web](#).

los hogares es asimétrica, por lo que sólo unos pocos hogares, especialmente los más ricos, son los que invierten en ciertos activos. Por otro lado, el componente panel indica que se vuelve a entrevistar a hogares que participaron en ediciones anteriores. Esto permite monitorearlos durante períodos de hasta diez años, y observar los cambios en las variables de interés de la encuesta. El número máximo de ediciones en las que un hogar puede participar en la EFF es de cuatro olas consecutivas. Si cesa su participación antes de completar sus cuatro ediciones, se descarta de la muestra y no vuelve a ser contactado en olas posteriores. Finalmente, para poder combinar ambas características y mantener la representatividad de la muestra para cada edición, en cada ola nueva se incluye una muestra de refresco con hogares nuevos.

En este proyecto vamos a utilizar datos que provienen tanto de las respuestas de los hogares al cuestionario de la EFF, como del paradata recopilado durante la elaboración de la encuesta. Para esto, es importante comprender bien cómo es el proceso de producción de la EFF y los datos que éste genera más allá de las respuestas de los hogares. Éste proceso puede verse de manera visual y resumida en la figura 4.1. La producción de la EFF se divide en dos grandes fases: Campo y Post-Campo. Durante la fase de Campo se contacta con los hogares, se realiza la entrevista personal, se procesan los datos y se realiza parte de la revisión de los mismos. En el Post-campo se termina el proceso de revisión, se revisa el grado de no-respuesta que hay en los datos, y se procede a la imputación de todas esas variables. Tras la imputación, se obtienen los datos finales. A continuación vamos a comentar con más detalle y en orden cronológico algunos aspectos de estas fases.

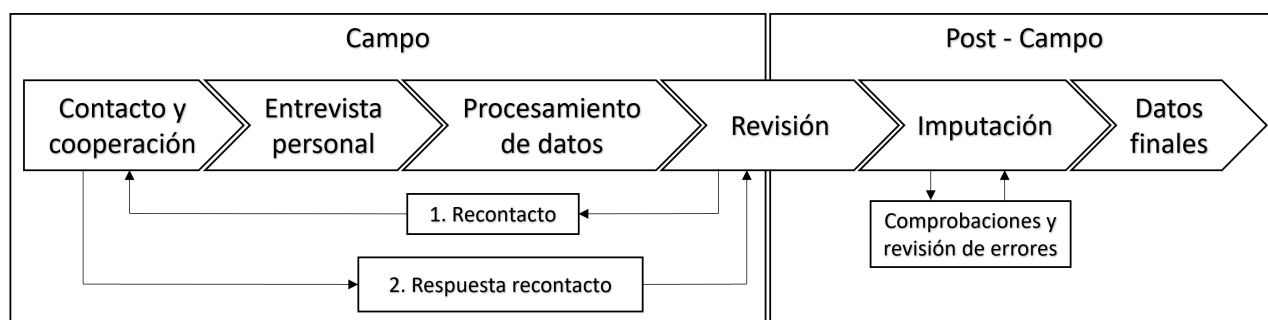


Figura 4.1: Fases de la creación de datos de la EFF

Para establecer el contacto, los entrevistadores realizan visitas personales a los hogares en sus domicilios. Esto puede requerir varios intentos, ya que es necesario que algún miembro del hogar esté físicamente en su hogar cuando tenga lugar la visita presencial. A continuación, se realiza una entrevista personal a la persona con más conocimiento de las finanzas del hogar, que se denomina Persona de Referencia o PR. La PR puede ser miembro del hogar, o un representante del mismo, denominado proxy, siempre y cuando sea la persona con mayor conocimiento de las

finanzas del hogar. También pueden participar en la entrevista otros miembros del hogar. Al principio de la entrevista, se pide a la PR su conformidad para grabar la entrevista en audio por motivos de calidad de los datos. Si se niega, la entrevista continúa, pero su revisión no dispondrá de audios. Los entrevistadores recogen las respuestas en un ordenador o una tablet (CAPI), y también anotan con comentarios de texto todos los detalles que consideren importantes para la revisión. La PR puede decidir no contestar a ciertas preguntas. Su valor se asignará como missing, y se imputará más adelante.

Adicionalmente, y sin la presencia de los hogares, los entrevistadores rellenan dos cuestionarios adicionales. Uno para recoger información sobre las características del vecindario y el tipo de edificio en el que vive el hogar, y otro con información sobre el desarrollo de la entrevista, como por ejemplo el nivel de interés mostrado por la PR o el nivel de entendimiento percibido de las preguntas.

Es necesario mencionar que algunos elementos de este procedimiento se modificaron durante la EFF2020, ya que el campo empezó en noviembre de 2020 y se vio afectado por la pandemia del Covid-19. Durante esa ola, los entrevistadores siguieron visitando personalmente a los hogares para conseguir su colaboración<sup>2</sup>, pero siempre respetando las medidas de distancia social. Las entrevistas se realizaron de manera telefónica asistida por una tablet (CATI). El resto del procedimiento se mantuvo.

Tras terminar la entrevista, todos los datos son procesados y enviados al equipo de revisión, que revisa individualmente todas las entrevistas y corrige los errores que puedan detectarse. Si se considera que hay información que se ha recogido erróneamente o se ha omitido, se recontacta con el hogar para corregir o recuperar esa información. Este recontacto se hace por teléfono. Tras este proceso de revisión, se analiza la proporción de preguntas sin responder de cada entrevista, y se eliminan las que no superen ciertos umbrales de calidad. Finalmente, todas las variables que contienen no-respuesta se imputan mediante técnicas de imputación múltiple<sup>3</sup>.

Para la elaboración de este proyecto se utiliza información recopilada durante las olas de la EFF2017, EFF2020 y EFF2022. La razón para elegir sólo estas ediciones es que son las únicas que cuentan con información detallada de parados y de las características de los entrevistadores. Para ediciones anteriores, esta información o no está disponible o contiene grandes errores de medida. A pesar de esto, es posible identificar el número de ediciones en las que ha participado cada hogar, por lo que es posible utilizar información que se remonta hasta la edición de la EFF2011.

A continuación se enumeran los ficheros de datos que se han utilizado durante este proyecto:

---

<sup>2</sup>Durante el principio del campo, algunos hogares panel fueron contactados sólo por teléfono, pero a las pocas semanas se decidió establecer la visita personal como el procedimiento estándar.

<sup>3</sup>Los métodos de imputación múltiple utilizados en la EFF pueden consultarse en [Barceló \(2006\)](#)

- Fichero de trabajo con respuestas de los hogares: Registro de hogares que contiene las respuestas al cuestionario justo antes de la imputación. Incluye las variables que contienen no-respuesta.
- Fichero de datos imputados: Registro de hogares que contiene las respuestas al cuestionario después de la imputación. Por simplicidad, sólo utilizaremos uno de los conjuntos de datos de la imputación múltiple.
- Fichero de contactos: Registro de hogares que contiene la información sobre el proceso de contacto con el hogar y sus resultados. Incluye el número de intentos de contacto, la fecha en que se produjo y el resultado de cada uno de ellos (aplazamientos, rechazos...).
- Fichero de recontactos: Registro de hogares que contiene información sobre el proceso de revisión y los recontactos.
- Censo de entrevistadores: Registro de entrevistadores que contiene la información disponible sobre los entrevistadores que han participado desde la EFF2014 a la EFF2022.
- Fichero de paradata: Registro de las pantallas del software CAPI o CATI. Para cada hogar, contiene información detallada de la interacción del entrevistador con el ordenador durante la entrevista. En concreto, contiene el flujo de pantallas que se visualizaron en el ordenador en cada entrevista y el tiempo que se estuvo en cada una.

## 4.2. Metodología

Tomando como referencia las estrategias de investigación propuestas en [Oates et al. \(2022\)](#), en este proyecto se sigue la estrategia de 'Caso de estudio', ya que se busca tener un conocimiento profundo sobre la participación de los hogares panel en el caso específico de la EFF, y buscar maneras de tratarlo. En concreto, vamos a intentar adaptar la implementación hecha por [Beste et al. \(2023\)](#) al caso de la EFF.

El objetivo es predecir si un hogar que ha participado en la ola  $t$  no volverá a participar en la siguiente ola  $t + 1$ . Para ello, se define una variable objetivo dicotómica Attrition que toma valor 0 si el hogar vuelve a participar, y valor 1 si no vuelve a participar:

$$Attrition = \begin{cases} 0 & \text{el hogar participa en la ola } t + 1 \\ 1 & \text{el hogar no participa en la ola } t + 1 \end{cases} \quad (4.1)$$

La metodología de este proyecto se fundamenta en cuatro pilares: La recopilación de los datos de los hogares y los entrevistadores que participaron en las ediciones de 2017, 2020 y 2022 de la

EFF, la selección de variables para los modelos de predicción a través de un análisis cualitativo y cuantitativo, el entrenamiento de varios modelos basados en métodos de machine learning con datos de la EFF2017 y la EFF2020 y la evaluación de su rendimiento para predecir la participación de los hogares panel en la EFF2022, y finalmente una valoración de los resultados obtenidos y qué pasos podrían darse en el futuro para seguir desarrollando este proyecto.

A continuación haremos una breve descripción sobre cómo ha sido el proceso de recopilación de datos. Y luego cerraremos este capítulo con la descripción de la estrategia de entrenamiento y validación de los modelos.

#### 4.2.1. Recopilación de datos

En la EFF hay, por un lado, información sobre los hogares repartida en varios ficheros para cada ola en la que han participado; y por otro, información sobre los entrevistadores en un censo de entrevistadores de la EFF. La unidad de análisis del estudio son hogares que ya han participado al menos una vez en alguna edición de la EFF, y que son elegibles para ser entrevistados en la siguiente edición. Por tanto, para cada ola se selecciona sólo a los hogares que han participado como máximo en tres ediciones. Los datos de esa ola se utilizarán como variables explicativas en los modelos, y como variable objetivo se extraerá el estatus de participación de ese mismo hogar del fichero de contactos de la ola siguiente. La vinculación de la información entre ambas ediciones se realiza a través de un identificador común que tienen los hogares en los ficheros de datos de ambas olas.

Por otro lado, la vinculación con la información del censo de entrevistadores se realiza a través de un identificador único que tiene cada entrevistador y que permanece inalterable a lo largo de las olas. Para cada hogar se conoce el identificador de la persona que realizó la entrevista en cada una de las olas en las que participó.

#### 4.2.2. Algoritmos de Machine Learning: Entrenamiento, validación y test

La estrategia de entrenamiento y metodología aplicada en este proyecto se basa en la aplicada en [Beste et al. \(2023\)](#). Para la evaluación de los modelos de panel attrition, el procedimiento habitual consiste en la elección de una regresión logística (Logit) como modelo base, y se realiza una comparación del rendimiento de cada modelo de machine learning con este modelo base ([Lepkowski et al. \(2002\)](#), [Kern et al. \(2021\)](#), [Beste et al. \(2023\)](#)). En la implementación de este documento se toma como modelo base una regresión logística de efectos primarios (es decir, sin interacciones entre variables).

Los modelos de machine learning que se evaluarán son un árbol de decisión (CART), un



Random Forest (RF), un eXtreme gradient boosting (XGBOOST) y un Naive Bayes (NB). La elección de estos modelos se debe a los buenos resultados que han mostrado para predecir panel attrition en otros estudios ([Kern et al. \(2019\)](#), [Kern et al. \(2021\)](#), [Beste et al. \(2023\)](#)), y a que son fácilmente interpretables en comparación con otros algoritmos de machine learning, como pueden ser las máquinas de soporte vectorial(SVM) o las redes neuronales. Esto facilita que sean utilizados en diseños adaptativos.

Con respecto al proceso de entrenamiento y evaluación, a la hora de elegir los conjuntos de entrenamiento, validación y test se ha tenido en cuenta el uso que se querría dar al algoritmo en la práctica. Siguiendo la filosofía de los diseños adaptativos, se quiere predecir con antelación qué hogares tienen más probabilidad de abandonar el estudio en la ola siguiente, y utilizar esa información para diseñar algún mecanismo dirigido a esos hogares. Eso implica establecer un criterio de separación temporal entre los conjuntos de entrenamiento y test, y asegurar que en los predictores de los modelos no se incluye información desconocida de antemano. En ese sentido, siguiendo la idea propuesta en [Beste et al. \(2023\)](#), utilizaremos exclusivamente información disponible en una ola concreta para predecir el resultado de la participación en la siguiente. Por esa razón, para el entrenamiento y la validación de los modelos, se utiliza la información recopilada en la EFF2017 para predecir la participación en EFF2020, y para el ejercicio del test se utiliza la información recopilada en la EFF2020 para predecir la participación en la EFF2022.

Antes de explicar cómo se procede con el tuning de hiperparámetros, es necesario hablar sobre la distribución de la variable Attrition. En el cuadro 3.1 puede observarse que, en las olas de 2020 y 2022, la mayoría de los hogares panel volvieron a participar. Es importante destacar que, aunque en la EFF2020 parezca que existe sólo un ligero desbalance entre las observaciones de cada clase de Attrition, durante los primeros entrenamientos y evaluaciones de los modelos se comprobó que las predicciones se asignaban masivamente a la clase mayoritaria (*Attrition* = 0). Para evitar esto, y dado el tamaño limitado de la muestra, previo al entrenamiento se procede al balanceo de la clase Attrition realizando un sobre-muestreo de dicha clase usando el algoritmo SMOTE (Syntetic Monirity Over-sampling Technique, [Chawla et al. \(2002\)](#)).

Attrition	EFF2020	EFF2022
Participa	3831	3974
No participa	2107	1531

Cuadro 4.1: *Distribución de Attrition en EFF2020 y EFF2022*

Para optimizar el tuning de los hiperparámetros de los modelos, se sigue una estrategia de k-fold cross-validation para los conjuntos de entrenamiento y validación, con k=5. Con respecto a los valores específicos de los hiperparámetros, como no existe certeza sobre cuáles son los

valores más adecuados, se considera un rango bastante amplio para cada hiperparámetro de todos los algoritmos. Para los modelos Logit y CART se utiliza un algoritmo de búsqueda de red (grid search) para probar cada combinación de hiperparámetros. Con respecto a los algoritmos RF y XGBOOST, dados los elevados tiempos y costes de ejecución de cada uno, se utiliza un algoritmo de búsqueda aleatoria (random search) con 2500 iteraciones. Aunque no se prueben todos los valores de los hiperparámetros, se ha comprobado que la búsqueda aleatoria es una alternativa válida en contextos de alta dimensionalidad de hiperparámetros ([Bergstra and Bengio \(2012\)](#)). Finalmente, a la hora de seleccionar las mejores combinaciones de modelos, la métrica que se busca maximizar es la ROC AUC.



## Capítulo 5

# Resultados y posibles próximos pasos

Este capítulo está compuesto por dos secciones. En la primera sección se muestran varios resultados y patrones observados durante el análisis exploratorio de los datos de la EFF, y cómo algunos de ellos han influido en la toma de decisiones con respecto a las variables que finalmente han entrado en los modelos. Finalmente, en la segunda sección se comentan los resultados de la evaluación de los diferentes modelos de machine learning utilizados durante este estudio.

Antes de empezar con los resultados, es importante destacar la gran cantidad de variables e información que se recoge en la EFF. En el Cuadro 4.1 se recoge el número de registros y variables disponibles en cada uno de los ficheros utilizados, después de seleccionar sólo a los hogares elegibles para la EFF2020 y a los entrevistadores que participaron en la EFF2017. El número de hogares elegibles para la EFF2020 es de 5.938, el número de entrevistadores es 69 y el número de pantallas que se vieron en el CAPI durante todas las entrevistas es 2.578.926<sup>1</sup>.

Nombre del fichero	Registros	Variables
Fichero de trabajo	5.938	6.103
Fichero de datos imputados	5.938	659
Fichero de recontactos	5.938	4
Fichero de contactos	5.938	636
Censo de entrevistadores	69	56
Fichero paradata	2.578.926	13

Cuadro 5.1: *Observaciones y variables de los ficheros de la EFF2017*

Con respecto al número de variables, muchas de ellas en su estado original no son informativas y necesitan ser combinadas con otras para poder obtener información interpretable. Esto reduce considerablemente el número de variables que hay que manejar. Otro grupo importante

---

<sup>1</sup>Los registros del fichero de paradata son las pantallas del ordenador con las que interactúa el entrevistador durante la entrevista. Suele haber una pregunta por pantalla, y cuando se responde, se pasa a la siguiente pantalla. Cuantas más preguntas, más pantallas se ven.

de variables tienen pocos registros porque se refieren por ejemplo inversiones en activos financieros que sólo tienen unos pocos hogares, y muchas de ellas se descartan, o se adaptan para indicar tenencia de dichos activos. Finalmente, todas las variables que aparecen en el fichero de datos imputados también aparecen en el fichero de trabajo<sup>2</sup>. Del fichero de datos imputados se usan las variables con los valores missing imputados, y el fichero de trabajo se utiliza para obtener indicadores de calidad de los datos, indicadores de no-respuesta y otras variables de interés que no aparecen en el fichero de datos imputados.

## 5.1. Análisis exploratorio

El análisis exploratorio de las variables es fundamental para poder ver cómo son las distribuciones de las variables que se van a utilizar, las relaciones que hay entre ellas, y especialmente la que hay con la variable de interés. También ayuda a detectar posibles errores o particularidades que puedan afectar a los modelos de predicción. Dada la enorme cantidad de variables que se han analizado en este proyecto, en este apartado sólo se presentan algunos resultados que han llamado bastante la atención y se considera que merecen la pena ser nombrados porque ofrecen conocimiento sobre la participación de los panelistas en futuras olas y que podría ser útil en la práctica.

En la figura 4.1 presenta cuatro gráficos de mekko que presentan las proporciones de hogares panel que vuelven a participar (región azul) frente a los que no vuelven a hacerlo (región roja) para las variables de número de olas en las que se ha participado (arriba a la izquierda), el nivel de recelo después de la entrevista (arriba a la derecha), el estado de salud reportado por la PR (abajo a la izquierda) y la edad de la PR (abajo a la derecha).

En la figura puede observarse que la proporción de abandono de los hogares que han participado sólo una ola es mayor que la de los que han participado en más de una. Este resultado sugiere que los hogares que llevan menos tiempo en la encuesta podrían ser más complicados de retener, y que podría ser interesante hacer un análisis enfocado en este tipo de hogares.

El segundo resultado a destacar es el del recelo mostrado por el hogar tras la entrevista. Recordemos que en la EFF se hacen preguntas que pueden considerarse delicadas, como por ejemplo si se poseen joyas, o cuál es el saldo que el hogar tiene en su cuenta corriente. Puede ser comprensible que se genere recelo. Este se mide con tres niveles: nada receloso, algo receloso y muy receloso. En el gráfico se observa que la proporción de hogares que abandonan la EFF aumenta con el nivel de recelo. Este resultado es bastante obvio, pero es muy importante

---

<sup>2</sup>El fichero de trabajo es el que recoge todas las correcciones de la revisión y es el que se utiliza para imputar los valores que están missing porque los hogares no los han declarado. Por tanto, la diferencia en las variables que ambos comparten es que en el fichero de datos imputados esas variables están completas porque sus valores missing se han imputado.

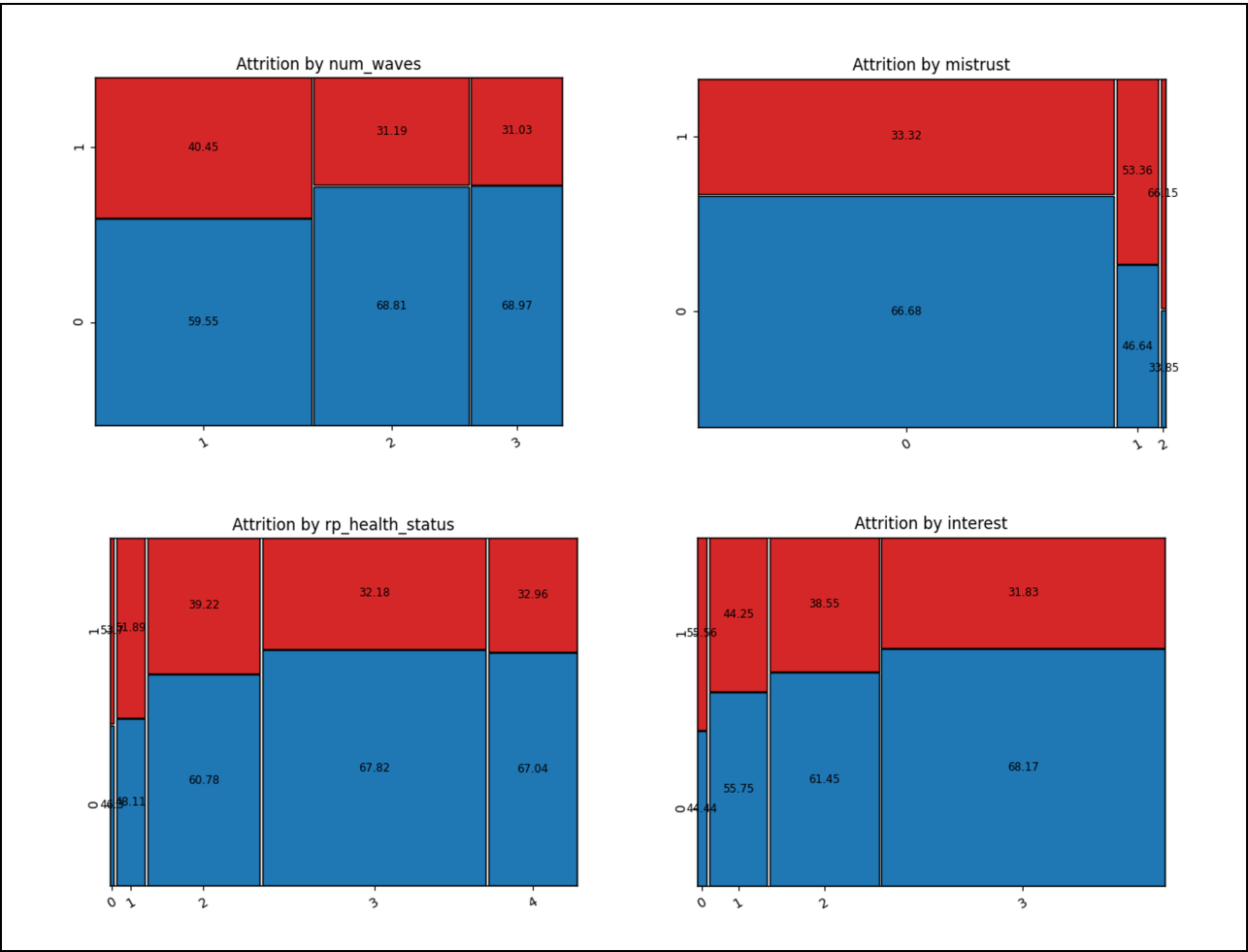


Figura 5.1: Categorical

resaltarlo porque a un entrevistador podría resultarle muy útil saber si el hogar con el que está a punto de contactar se mostró muy receloso después de hacer la entrevista en la ola anterior, y en su estrategia de contacto y de ganar la cooperación podría hacer más énfasis en los aspectos de la encuesta que causan más recelo.

Finalmente, los dos gráficos de la parte de inferior hacen referencia a características de la PR. En el de la izquierda se observa que la proporción de hogares que abandona la encuesta es mayor cuando menor es el nivel de salud que reporta la PR. De nuevo, esto puede parecer algo obvio porque alguien que tiene peor salud seguramente o quiera gastar su tiempo en realizar una encuesta. Pero, de nuevo, la estrategia de contacto y de ganar la cooperación de un entrevistador puede adaptarse si ya saben que van a ver a alguien que seguramente tenga mala salud, ya sea ofreciendo hacer la entrevista en varias sesiones, o en un ambiente en el que la PR se muestre más cómoda.

Finalmente, el último gráfico muestra el nivel de interés que mostró la PR durante la entrevista durante la edición anterior. Se observa que la proporción de hogares que vuelven a participar es mayor cuanto mayor es el interés que mostraron durante la ola anterior. De nuevo, se trata de un resultado obvio. Pero, al igual que con el nivel de recelo, este podría ser un dato que podría ser muy útil para el entrevistador cuando esté preparando su estrategia de contacto con el hogar.

## 5.2. Evaluación de modelos

Para la selección de las variables que han entrenado los modelos se han seguido X criterios. En primer lugar, se han seleccionado variables que aparecen en las implementaciones de [Kern et al. \(2021\)](#) y [Beste et al. \(2023\)](#) en sus modelos de machine learning. En segundo lugar, se han seleccionado variables sobre factores mencionados en la revisión de [Lynn \(2018\)](#) y que han mostrado potencial para predecir la no participación de hogares panel, como por ejemplo la duración o información sobre los contactos con los hogares. En tercer lugar, se han seleccionado variables específicas que se recogen en la EFF, como la tenencia de activos, deudas o el recelo de los hogares percibido por los entrevistadores. Finalmente, tras seleccionar todas las variables anteriores, se han realizado diversos tests gráficos y estadísticos para detectar variables con correlaciones o dependencias muy altas, y se han descartado. La selección final de variables es la que aparece en el cuadro 4.2.

El cuadro 4.3 contiene los resultados del rendimiento de los modelos entrenados sobre el conjunto de test, es decir, los resultados para predecir la no participación de los panelistas en la EFF2022. Las métricas de evaluación que se utilizan son Accuracy, Precision, Recall, F1 y ROC AUC. La métrica de referencia que utilizamos para la evaluación es la ROC AUC, que

Fuente de información	Variables
Características del hogar	Número de adultos que trabajan, Número de adultos jubilados, Propietario vivienda principal, Tamaño del hogar, Tiene otras propiedades, Tiene joyas, Posee negocios, Posee cuentas para pagos, Posee acciones que cotizan, Posee acciones que no cotizan, Posee renta fija, Posee fondos de inversión, Posee cuentas para no pagos, Posee planes de pensiones, Les deben dinero, Poseen vehículos, Percentil de renta, Percentil de riqueza Bruta, Pareja vive en el hogar, Poseen otros activos financieros, Tienen deuda, Tienen ingresos de activos, Hijos viven en el hogar, Nivel de satisfacción con la vida
Características PR	Nivel de satisfacción con la vida, PR es panel, PR - nivel educativo, PR- Edad, PR - Casada, PR - Viuda, PR - Sexo, Situación laboral - Asalariado, Situación laboral - Jubilado, Situación laboral - Inactivo, PR - Estado de salud
Valoraciones FI	Recelo tras la entrevista, Tipo de edificio - Unifamiliar, Barreras - portero automático, Barreras - No hay barreras, Entendimiento de las preguntas PR, Interés PR, Razones colaborar - Interesado en estos estudios, Razones para colaborar - Lo lleva el BdE, Razones para colaborar - Relevancia de la encuesta, Razones - Favor al entrevistador, Razones - Dar su opinión, Razones - Otras
Paradata	PR consiente grabar la entrevista, Hogar con recontacto exitoso, Número de olas en las que se ha participado, Tamaño del municipio, Entrevista con proxy, Número de miembros del hogar que participan, Ratio estricto, Ratio amplio, Número de variables con valores missing, Duración de la entrevista

Cuadro 5.2: Selección de variables para entrenar los modelos de predicción



se encuentra en la parte derecha del cuadro. Esta métrica mide el rendimiento entre la tasa de falsos positivos y falsos negativos. Toma valores de 0.5 a 1, con 1 siendo un predictor perfecto, y 0.5 el que se obtendría con una estimación realizada de manera aleatoria.

Modelo	Accuracy	Precision	Recall	F1	ROC AUC
Logit	0,6556	0,3749	0,3573	0,3659	0,5959
CART	0,6434	0,3338	0,2835	0,3066	0,5521
Random Forest	0,6489	0,3529	0,3148	0,3328	0,5821
XGBooster	0,6718	0,3772	0,2769	0,3194	0,5911
Naive Bayes	0,6254	0,3504	0,4063	0,3763	0,5798

Cuadro 5.3: Métricas de evaluación de los modelos de predicción en el conjunto de test

El modelo de referencia de este estudio, el Logit, presenta una ROC AUC de 0,5959. Se considera que un valor inferior a 0,6 es un resultado malo, por lo que el modelo Logit no es un buen predictor. Con respecto a los otros modelos, observamos que el CART y el Naïve Bayes presentan valores más bajos que el Logit, de 0,5521 y 0,5798 respectivamente. El Random Forest y el XGBooster, en cambio, presentan valores de 0,5821 y 0,5911 respectivamente, que son ligeramente superiores al del modelo Logit.

De estos resultados podemos ver que los modelos de Random Forest y XGBooster mejoran al Logit de referencia. Y entre estos dos, el Random Forest presenta valores un poco más altos en Accuracy y en Precision, y el XGBooster es mejor en Recall y F1. Aun así, es necesario recalcar que el valor de ROC AUC sigue estando por debajo de 0,6, por lo que, aunque se mejore el rendimiento con respecto al modelo Logit, los resultados de los test son malos.

Estos resultados indican que hay modelos de machine learning que mejoran en rendimiento de la predicción con respecto a un Logit. Sin embargo, las métricas señalan que el rendimiento del mejor modelo es regular. La pregunta que hay que preguntarse entonces es, ¿por qué no está funcionando bien la predicción?

Una de las causas puede ser el overfitting, es decir, que los modelos no son capaces de generalizar bien porque se han adaptado demasiado bien a los datos con los que han sido entrenados, y por tanto su rendimiento no es bueno cuando se les presenta con datos nuevos. una manera de comprobar esto es haciendo el ejercicio predicción sobre los datos de entrenamiento con estos mismos modelos, y ver cómo son esos resultados. Estos resultados pueden verse en el cuadro 4.4.

Como era de esperar, los resultados de las predicciones de los modelos sobre los datos de entrenamiento son mejores que los que se observan con respecto a los datos test. En este caso, el modelo que presenta mejor rendimiento es el Random Forest, con una métrica de ROC AUC de 0,6726, ligeramente superior a la del XGBooster. El resto de métricas del Random Forest también son ligeramente mejores que las del XGBooster. Sin embargo, esta métrica de ROC

Modelo	Accuracy	Precision	Recall	F1	ROC AUC
Logit	0,6389	0,4905	0,4518	0,4704	0,6517
CART	0,6126	0,4594	0,5178	0,4868	0,6188
Random Forest	0,6596	0,5223	0,4770	0,4986	0,6726
XGBooster	0,6544	0,5144	0,4675	0,4898	0,6722
Naive Bayes	0,6251	0,4741	0,5178	0,4950	0,6435

Cuadro 5.4: Métricas de evaluación de los modelos de predicción en el conjunto de entrenamiento

AUC está entre 0,6 y 0,75, que lo clasifica como un test regular.

Continuando con el Random Forest, una ventaja de los modelos basados en árboles con respecto a otros modelos es que pueden ser interpretados. En el caso del Radom Forest, es posible consultar qué variables han tenido más peso a la hora de clasificar la participación de los hogares. Aunque los resultados de la predicción del training no sean buenos, merece la pena echarle un ojo por los patrones que haya podido detectar. El cuadro 4.5 presenta las 10 variables con más importancia en el modelo de Random Forest.

Variable	Importancia
PR es panel	0,0822
Interés PR	0,0790
Razones para colaborar - Relevancia de la encuesta	0,0605
Ratio amplio	0,0594
Poseen vehículos	0,0539
Tienen deuda	0,0496
Situación laboral - Asalariado	0,0475
Razones colaborar - Interesado en estos estudios	0,0380
Posee planes de pensiones	0,0353
PR consiente grabar la entrevista	0,0348

Cuadro 5.5: Selección de variables para entrenar los modelos de predicción

La importancia de una variable en el Random Forest mide el peso relativo que ha tenido una variable concreta a la hora de crear las ramificaciones de los diferentes árboles de decisiones que va generando el Random Forest durante su entrenamiento. En este caso de la participación en la EFF2020, las tres variables que más importancia tuvieron fueron que la PR fuera panel, que la PR mostrase interés durante la entrevista, y que el entrevistador indicase que el hogar colaboró por la relevancia de la encuesta. También es interesante destacar la variable ratio amplio, que es un indicador de la no-respuesta del hogar a la hora de facilitar cantidades monetarias<sup>3</sup>. Finalmente, otra variable que es importante para hacer la clasificación es si la PR consintió que

<sup>3</sup>El ratio amplio se define como el cociente entre el número de preguntas en euros respondidas por el hogar, ya fuera como valor puntual o como intervalos, sobre el número total de preguntas planteadas. Cuando mayor es su valor, más información ha facilitado el hogar.

se grabase la entrevista.

## Capítulo 6

# Conclusiones y futuras líneas de trabajo

En este Trabajo de Final de Master se ha intentado aplicar la implementación de [Beste et al. \(2023\)](#) para predecir la no participación de hogares panel en encuestas longitudinales en olas futuras al caso de la Encuesta Financiera de las Familias. Para ello, se ha usado un modelo de Regresión Logística como referencia, y se han entrenado modelos CART, Random Forest, XGBooster y Naïve Bayes. El conjunto de predictores está formado por variables que potencialmente pueden explicar la participación de un hogar en la edición siguiente de la encuesta, como por ejemplo el interés mostrado por la persona que respondió a la encuesta, el nivel de recelo que mostró durante la entrevista, su nivel de salud o si consintió que la entrevista fuese grabada.

Los modelos de Random Forest y XGBooster presentaron mejores rendimientos que el modelo de referencia Logit en todas las métricas de evaluación consideradas, pero el valor de ROC AUC del mejor de los modelos no superó el valor de 0.6, lo cual lo clasifica como un predictor malo. Para intentar aprender sobre cómo se ha hecho la predicción, se hizo el ejercicio de usar los mismos modelos para hacer la predicción con el conjunto de entrenamiento. El modelo que funcionó mejor fue el Random Forest, pero su ROC AUC no superó el valor de 0.7, que lo clasifica como un predictor regular. Al observar las variables que más importancia tuvieron durante el entrenamiento del Random Forest destabacan que la persona que contestó a la entrevista a ya formase parte del hogar desde al menos dos ediciones antes, el valor del interés que mostró también era importante, y que el entrevistador considerase que el hogar participó por la relevancia de la encuesta.

A partir de los resultados que se han visto en este proyecto, se plantean las siguientes reflexiones y los posibles pasos que se podrían dar en el futuro para este proyecto:

1. La exploración de los datos sugiere que las variables seleccionadas para entrenar los mo-

delos de predicción guardan cierta relación con la variable de attrition. Pero no son suficientes para predecir bien si un hogar panel dejará de participar en la ola siguiente. Tal y como hemos comentado al principio de este capítulo, la EFF genera una gran cantidad de variables, y se hizo un filtrado inicial basado en las implementaciones de Beste et al. (2023) y Kern et al. (2021). Es posible que haya **variables que no se han incluido, y que tengan poder para predecir el attrition**. Una vía de trabajo para el futuro es **volver a revisar toda la información disponible y plantear una nueva selección de variables**. El uso de métodos de machine learning para selección de variables es una opción que se podría implementar.

2. Las olas de la EFF seleccionadas para el estudio son las de los años 2017, 2020 y 2022 porque son las que tienen mayor cantidad de datos y también los de mayor calidad. Todos los modelos entrenados buscan utilizar datos de lo que había en 2017 para predecir algo que iba a ocurrir en 2020, que es el año en el que tuvo lugar la crisis del **Covid-19**. Y posteriormente, en el test, se utiliza información recogida durante 2020 para predecir lo ocurrido en el año 2022, cuando el mundo estaba ya superando la crisis del Covid-19. El año 2020 fue un año especial en la EFF porque seguramente mucha gente fuera más reticente a participar por el covid, lo cual afectaría a la variable target del entrenamiento). También se tuvo que cambiar la metodología de la entrevista, pasando de entrevista personal a entrevista telefónica, que afecta a la calidad de los datos (Lynn (2018)), y por tanto a los datos de los predictores usados en el test. Ante esta problemática, se plantean dos posibles alternativas:
  - a) Las olas de **EFF2002 a EFF2017 son homogéneas** en lo que a metodología se refiere. Todas las entrevistas fueron personales, y no hubo una crisis como la del Covid-19. Una alternativa interesante sería crear **nuevos modelos de predicción**, pero utilizando sólo la información disponible en esas olas. Serían **menos variables** (las respuestas de los hogares y la información recogida por los entrevistadores), pero podría ser que el rendimiento fuese mejor.
  - b) La EFF va a continuar realizándose en los próximos años, y con una frecuencia bienal. Se puede volver a **plantear esta misma metodología dentro unos años, utilizando sólo ediciones completadas después de 2020**, y con el beneficio de recoger toda la información que se ha estado recogiendo en las últimas ediciones y que no está disponible para antes de 2017.
3. Los resultados de las predicciones sobre los datos de entrenamiento también podría explicarse por la **existencia de información no observable en el momento de hacer la predicción**, y que además tenga más peso para determinar el resultado de la

**participación que la información recabada durante la ola anterior.** El Covid-19 es un buen ejemplo de algo que no se puede prever, pero otra cosa que también se desconoce es qué entrevistadores habrá durante la siguiente edición. En todas las ediciones hay entrevistadores nuevos, y algunos funcionan muy bien, y otros no. Tal y como comentan [Lynn \(2018\)](#) y [Groves \(2006\)](#), el papel del entrevistador es importante para la colaboración de los hogares y la calidad de los datos. Esto puede investigarse identificando a los hogares para los que no se han hecho buenas predicciones, y analizar por un lado cómo son las características que tienen en la ola anterior, y por otro la información que se tenga sobre la ola que se intenta predecir, y comprobar si las variables que tienen más peso para explicar la participación son las de la ola corriente o las de la ola anterior.

4. Una opción que siempre hay que considerar es un **cambio de enfoque**. Hay dos alternativas que son interesantes:
  - a) En la exploración de los datos vimos que los hogares que han participado sólo en una edición muestran más proporción de abandonos que los que han participado más de dos años. Una posible explicación de esto es que los hogares que han participado más de una vez están más comprometidos con el estudio, y seguramente merezca la pena **enfocar el análisis en la predicción de la participación de los paneles en su segunda ola** en vez de hacer la predicción para todos los hogares.
  - b) En vez de predecir un resultado binario, de participar o no participar, se puede plantear hacer un ejercicio de **análisis de supervivencia** (survival analysis), e intentar predecir el número de ediciones en las que participará un hogar de la EFF antes de abandonar el estudio. Esta información está disponible y se podría utilizar información de todas las olas de la EFF.



# Bibliografía

- Alvargonzález, P., Pidkuyko, M., and Villanueva, E. (2020). La situación financiera de los trabajadores más afectados por la pandemia: un análisis a partir de la encuesta financiera de las familias. *Boletín Económico/Banco de España*, 3/2020.
- Barceló, C. (2006). Imputation of the 2002 wave of the spanish survey of household finances (eff). *Banco de Espana Research Paper No. OP-0603*.
- Barceló, C., Crespo, L., Garcia, S., Gento, C., Gómez, M., and de Quinto, A. (2021). The spanish survey of household finances (eff): Description and methods of the 2017 wave. *Banco de Espana Occasional Paper*, 2033.
- Bergstra, J. and Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of machine learning research*, 13, No. 2.
- Beste, J., Frodermann, C., Trappmann, M., and Unger, S. (2023). Case prioritization in a panel survey based on predicting hard to survey households by machine learning algorithms: An experimental study. In *Survey Research Methods*, volume 17, No. 3, pages 243–268.
- Bover, O. (2004). Encuesta financiera de las familias españolas (eff): Descripción y métodos de la encuesta de 2002. *Banco de España, Documentos Ocasionales*, 0409.
- Bover, O. (2008). The spanish survey of household finances (eff): description and methods of the 2005 wave. *Banco de España, Documentos Ocasionales*, 0803.
- Bover, O. (2011). The spanish survey of household finances (eff): description and methods of the 2008 wave. *Banco de España Occasional Paper*, 1103.
- Bover, O., Coronado, E., and Velilla, P. (2014). The spanish survey of household finances (eff): description and methods of the 2011 wave. *Banco de Espana occasional Paper*, 1407.
- Bover, O., Crespo, L., Gento, C., and Moreno, I. (2018). The spanish survey of household finances (eff): description and methods of the 2014 wave. *Banco de España Ocassional Paper*, 1804.



- Buskirk, T. D., Kirchner, A., Eck, A., and Signorino, C. S. (2018). An introduction to machine learning methods for survey researchers. *Survey Practice*, 11(1).
- Calderwood, L., Cleary, A., Flore, G., and Wiggins, R. (2012). Using response propensity models to inform fieldwork practice on the fifth wave of the millenium cohort study. In *International Panel Survey Methods Workshop*, pages 4–5.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P. (2002). Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357.
- Couper, M. P. and Ofstedal, M. B. (2009). Keeping in contact with mobile sample members. *Methodology of longitudinal surveys*, pages 183–203.
- Crespo, L., El Amrani, N., Gento, C. L., and Villanueva, E. (2023). Heterogeneidad en el uso de los medios de pago y la banca online: un análisis a partir de la encuesta financiera de las familias (2002-2020). *Documentos Ocasionales/Banco de España*, 2308.
- De Leeuw, E. D. et al. (2005). To mix or not to mix data collection modes in surveys. *Journal of official statistics*, 21(5):233–255.
- Early, K., Mankoff, J., and Fienberg, S. E. (2017). Dynamic question ordering in online surveys. *Journal of Official Statistics*, 33(3):625–657.
- Groves, R. M. (2006). Nonresponse rates and nonresponse bias in household surveys. *International Journal of Public Opinion Quarterly*, 70(5):646–675.
- Groves, R. M., Cialdini, R. B., and Couper, M. P. (1992). Understanding the decision to participate in a survey. *Public opinion quarterly*, 56(4):475–495.
- Groves, R. M., Fowler Jr, F. J., Couper, M. P., Lepkowski, J. M., Singer, E., and Tourangeau, R. (2009). *Survey methodology*. John Wiley & Sons.
- Groves, R. M. and Heeringa, S. G. (2006). Responsive design for household surveys: tools for actively controlling survey errors and costs. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 169(3):439–457.
- Gómez, M. and Villanueva, E. (2022). El efecto de los planes de pensiones de empresa sobre el ahorro privado de los hogares. *Boletín Económico/Banco de España*, 2/2022.
- Hirano, K., Imbens, G., Ridder, G., and Rebin, D. B. (1998). Combining panel data sets with attrition and refreshment samples.

- Jankowsky, K. and Schroeders, U. (2022). Validation and generalizability of machine learning prediction models on attrition in longitudinal studies. *International Journal of Behavioral Development*, 46(2):169–176.
- Kern, C., Klausch, T., and Kreuter, F. (2019). Tree-based machine learning methods for survey research. In *Survey research methods*, volume 13, No. 1, page 73. NIH Public Access.
- Kern, C., Weiß, B., and Kolb, J.-P. (2021). Predicting nonresponse in future waves of a probability-based mixed-mode panel with machine learning. *Journal of Survey Statistics and Methodology*, page smab009.
- Kreuter, F. and Müller, G. (2015). A note on improving process efficiency in panel surveys with paradata. *Field Methods*, 27(1):55–65.
- Lagorio, C. et al. (2016). Call and response: Modelling longitudinal contact and cooperation using wave 1 call records data. Technical report, Understanding Society at the Institute for Social and Economic Research.
- Laurie, H. and Lynn, P. (2009). The use of respondent incentives on longitudinal surveys. *Methodology of longitudinal surveys*, pages 205–233.
- Laurie, H., Smith, R. A., and Scott, L. (1999). Strategies for reducing nonresponse in a longitudinal panel survey. *Journal of Official Statistics*, 15(2):269–282.
- Lepkowski, J. M., Couper, M. P., et al. (2002). Nonresponse in the second wave of longitudinal household surveys. *Survey nonresponse*, pages 259–272.
- Liu, M. (2020). Using machine learning models to predict attrition in a survey panel. *Big data meets survey science: A collection of innovative methods*, pages 415–433.
- Lynn, P. (2009). Methods for longitudinal surveys. *Methodology of longitudinal surveys*, pages 1–19.
- Lynn, P. (2014). Targeted response inducement strategies on longitudinal surveys. *Improving Survey Methods*, pages 322–338.
- Lynn, P. (2017). From standardised to targeted survey procedures for tackling non-response and attrition. In *Survey Research Methods*, volume 11, No. 1, pages 93–103.
- Lynn, P. (2018). Tackling panel attrition. *The Palgrave handbook of survey research*, pages 143–153.

- Lynn, P., Kaminska, O., and Goldstein, H. (2014). Panel attrition: how important is interviewer continuity? *Journal of Official Statistics*, 30(3):443–457.
- Mcgonagle, K. A., Sastry, N., and Freedman, V. A. (2022). The effects of a targeted “early bird” incentive strategy on response rates, fieldwork effort, and costs in a national panel study. *Journal of Survey Statistics and Methodology*, page smab042.
- Nicoletti, C. and Peracchi, F. (2005). Survey response and survey characteristics: microlevel evidence from the european community household panel. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 168(4):763–781.
- Oates, B. J., Griffiths, M., and McLean, R. (2022). *Researching information systems and computing*. Sage.
- O’Brien, E. M., Black, M. C., Carley-Baxter, L. R., and Simon, T. R. (2006). Sensitive topics, survey nonresponse, and considerations for interviewer training. *American journal of preventive medicine*, 31(5):419–426.
- Rubin, D. B. (1987). *Multiple imputation for nonresponse in surveys*. New York: Wiley.
- Schouten, B., Peytchev, A., and Wagner, J. (2017). *Adaptive survey design*. CRC Press.
- Tourangeau, R., Michael Brick, J., Lohr, S., and Li, J. (2017). Adaptive and responsive survey designs: a review and assessment. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 180(1):203–223.
- Wagner, J. R. (2008). *Adaptive survey design to reduce nonresponse bias*. PhD thesis, University of Michigan.
- Watson, N. and Wooden, M. (2009). Identifying factors affecting longitudinal survey response. *Methodology of longitudinal surveys*, pages 157–181.