



**¿Sabe cuánto tiene ahorrado en su cuenta de
AFP? Caracterizando individuos con algoritmos
de Explainable Artificial Intelligence.**

**TESIS PARA OPTAR AL GRADO DE
MAGÍSTER EN FINANZAS**

**Alumno: Tomás Orellana Tapia
Profesor Guía: David Diaz Solis**

Santiago, Enero 2023



Abstract

El objetivo de este trabajo consiste en modelar si los individuos saben el monto que tienen ahorrado en su cuenta individual de jubilación (sistema AFP), lo que podría afectar su bienestar al momento de la vejez (bajos montos implicarían una menor pensión). Para realizar esto, utilizamos la Encuesta de Protección Social (EPS) del año 2015 en la cual realizan la siguiente pregunta a los individuos, *¿sabe usted cuánto hay acumulado en su cuenta individual?*. Como herramienta predictiva, y debido al alto grado de acierto alcanzado en tareas similares, este trabajo ocupa como herramienta algoritmos supervisados de Machine Learning (*ML*). Sin embargo, en la literatura se presenta un trade-off entre el alto poder predictivo alcanzado por estos algoritmos y su falta de transparencia de los patrones que explican la predicción (*Black-box*). Dado lo anterior, en este trabajo se utilizan los últimos métodos Post-Hoc de análisis disponibles para “abrir la caja negra”, de manera de entregar información respecto de los determinantes más comunes a nivel de la población, y también, cuantificar la importancia de las diferentes variables explicativas, y su sensibilidad a nivel individual.

Primeramente se utilizó *CatBoost* como el algoritmo predictivo, y luego se utilizaron técnicas de “*Explainable-AI*”, en particular *SHAP-values*, para calcular la importancia de cada variable predictiva utilizada en el modelo, a nivel poblacional e individual. Los resultados además se comparan con los obtenidos con un modelo estadístico tradicional de regresión logística.

Nuestros resultados indican que podemos clasificar correctamente al 74 % de los casos, y que las variables predictivas más importantes se relacionan con haber recibido información actualizada respecto del estado de su ahorro, y el nivel de conocimiento del funcionamiento del sistema de pensiones. Nuestros resultados y técnicas utilizadas pueden ser de gran utilidad para la generación de políticas públicas que mejoren el conocimiento respecto de los planes de retiro y jubilación.



Índice

1. Introducción	4
2. Revisión de literatura	6
3. Metodología, Pre-Proceso y Estadística Descriptiva	12
4. Resultados	16
4.1. Estadísticas de Clasificación.	16
4.2. Análisis del modelo.	20
4.3. Análisis de individuos.	24
4.4. Relevancia Económica	26
4.5. Discusión.	28
5. Conclusión	29
6. Anexos	34



1. Introducción

La teoría clásica del ciclo de vida del consumo señala que los individuos tendrán un comportamiento racional en el gasto que tienen a lo largo de su vida. Esto es, en un principio no tienen ingresos y se endeudan, luego obtienen ingresos y son capaces de pagar deuda y ahorrar para su jubilación, y finalmente gastan lo ahorrado durante su vida activa.

Esta teoría es difícil de sostener, ya que uno de los principales motivos por los cuales existe ahorro forzoso para la jubilación, es debido a la incapacidad de los individuos de postergar sus decisiones de consumo presente para obtener a cambio un mayor ingreso futuro.

Por otro lado, la literatura señala diversos beneficios de la educación financiera en el bienestar de los individuos. Por ejemplo, (Behrman y col., 2012) señalan que una mayor educación financiera aumenta la probabilidad de que los individuos contribuyan más al ahorro de su pensión, lo cual aumentaría su riqueza neta. Además, los autores encuentran que la educación financiera puede hacer contribuciones significativas al comportamiento financiero, mucho más allá de lo que puede contribuir la educación regular (*regular schooling*). Por parte de la econometría que realizan los autores, el usar MCO (Mínimos Cuadrados Ordinarios) para estimar el efecto de la educación financiera y la escolaridad sobre la riqueza puede conllevar a resultados engañosos debido a errores de medición y factores inobservables.

El trabajo de (Klapper y col., 2012) indica que la educación financiera está significativamente relacionada con una mayor participación en los mercados financieros formales y negativamente relacionado con el uso de préstamos informales. Además, una mayor educación financiera ayuda a los individuos a enfrentar de mejor manera shocks macroeconómicos y de ingresos inesperados. Los datos de panel que utilizan los autores corresponde a individuos de Rusia, en donde se puede desprender que solo el 41 % de los encuestados sabe cómo funciona el *interés compuesto* y solo el 46 % puede responder una simple pregunta sobre inflación.

Con respecto a individuos que están cercanos a la jubilación, se encuentra que estos tienen poco conocimiento de los planes de retiro, y esta falta de conocimiento afecta sus planes de jubilación (Clark y col., 2012). De hecho, conceptos erróneos sobre los planes de jubilación tienen una influencia en la edad esperada de jubilación que tienen los trabajadores. Estos errores de información son asimétricos, los trabajadores que creen que los beneficios de la jubilación se pueden obtener en una edad tardía, tienden a retirarse con una edad más avanzada.



Debido a lo anterior, este trabajo se enfoca en modelar y predecir si los individuos saben el monto que tienen ahorrado en su cuenta de jubilación. Para efectuar esto, utilizamos la Encuesta de Protección Social 2015 (EPS), en donde se le pregunta a los individuos si saben cuanto tienen acumulado en su cuenta individual del sistema de pensiones chileno (AFP).

Además, no solo entregamos la capacidad predictiva que tiene nuestro modelo, si no que también empleamos nuevas herramientas para poder explicar las predicciones que hace el modelo, esto es debido a que los resultados entregados son del tipo caja negra, es decir, no se sabe o es difícil de interpretar la forma en la que el modelo llegó al resultado entregado. En consecuencia, empleamos metodologías que nos permiten abrir esta caja negra y poder mostrar cuales fueron las variables más importantes en el modelo para poder realizar la predicción, e incluso podemos observar la importancia de las variables a nivel de cada individuo.

Dado esto, podemos explorar las variables de los individuos que tuvieron una respuesta negativa y tratar de revertir esta situación mediante mayor educación o información sobre los planes de jubilación. Esto puede ser el incentivo a crear un programa de política pública que mediante una pequeña entrega de información se pueda ayudar sustancialmente al individuo a conocer más sobre su situación financiera para la vejez.

Dentro de la base de datos, contamos con distintas variables socio-económicas (educación, salud, situación laboral) que utilizamos como input para nuestra predicción. Con estas variables y utilizando un algoritmo de *Árbol de decisión* llamado *Catboost*, predecimos nuestra variable de interés. Nuestro algoritmo pudo clasificar correctamente el 74% de las predicciones realizadas. Las variables más importantes en la predicción para nuestro modelo son: *en los últimos 12 meses, ¿ha recibido alguna cartola de su AFP?* y *¿sabe usted que porcentaje de su ingreso imponible le descuentan?*

Tal como se ha mencionado en el texto, utilizamos algoritmos de *Machine Learning* en este trabajo, por lo cual, cabe mencionar una breve explicación de lo que es este concepto. En primer lugar, debemos nombrar que Machine Learning corresponde a una rama de la *Inteligencia Artificial*. Este último campo ha tenido numerosos tipos de definiciones, pero una forma resumida de agrupar estos conceptos es la definición que entrega el profesor John McCarthy en su paper del año 2004: “La IA es la ciencia e ingeniería de hacer inteligentes a las máquinas, especialmente softwares informáticos inteligentes. Está relacionado con la tarea similar de usar computadoras para comprender la inteligencia humana, pero la IA no se tiene que confinar a los



métodos que son biológicamente observables” (McCarthy, 2004).

En lo referente a *Machine Learning*, este campo “se concentra en el uso de datos y algoritmos para poder imitar la forma en la que los humanos aprenden, mejorando gradualmente su precisión.” “...A través del uso de métodos estadísticos, los algoritmos son entrenados para hacer clasificaciones o predicciones, con el objetivo de descubrir perspectivas claves en proyectos con alto uso de datos, en donde posteriormente son utilizados para la toma de decisiones de aplicaciones y negocios.” (IBM, 2020).

Nuestro aporte en la literatura radica en que utilizamos nuevos algoritmos, especialmente para la industria financiera, con el objetivo de poder predecir cuanto sabe de su ahorro el individuo, dada la información que maneja. Además, no solo encontramos la predicción, sino que con estas nuevas herramientas podemos modificar las respuestas de los individuos con el objetivo de cambiar su output y obtengan una respuesta acorde a lo que necesitamos, es decir, que sepan cuanto hay acumulado en su cuenta de AFP.

Este trabajo está dividido en la siguiente forma, la sección 2 hace un análisis de la literatura relevante y relacionada con nuestro trabajo. En la sección 3 se presenta la metodología y Estadística descriptiva. La sección 4 muestra los resultados encontrados y se realiza un análisis tanto del modelo como de los individuos, además entregamos la relevancia económica y una discusión de nuestros resultados. Por último, la sección 5 entrega la conclusión del trabajo.

2. Revisión de literatura

En la introducción se dio una breve definición de lo que son los conceptos de *Inteligencia Artificial* y *Machine Learning*. si bien las definiciones dadas son contemporáneas, los cimientos de estos campos fueron dados en un principio por *Alan Turing* (considerado como padre de la ciencia de la computación y precursor de la informática moderna), quien en su artículo *Computing Machinery and Intelligence* publicado el año 1950 presenta lo que se conoce como el *Test de Turing*, el cual corresponde en determinar si una máquina es capaz de demostrar inteligencia humana. Para lograr esto, Turing propone lo siguiente: si una máquina logra entablar una conversación con un humano sin ser detectado como una máquina, entonces la máquina ha demostrado inteligencia humana (Turing, 2009). Por lo cual, el objetivo de estos algoritmos es imitar la complejidad de la inteligencia humana con el fin de solucionar u optimizar problemas de un alto nivel del uso de datos.



La aplicación de estos algoritmos se ha encontrado en diversos campos, por ejemplo, en el campo de la medicina se han utilizado técnicas del tipo *árbol de decisión* para poder predecir el cáncer de mama. Este es el enfoque del artículo de Venkatesan y Velmurugan, en donde utilizando como input información de los pacientes tal como *sexo, edad, altura, peso, historial pasado, diagnósticos médicos, hábitos alimenticios, entre otros*, son capaces de clasificar correctamente el 99% de los casos de cáncer de mama (Venkatesan y Velmurugan, 2015). Para ver más aplicaciones en la medicina, leer (Shailaja y col., 2018).

Podemos encontrar ejemplos también en el área de la agricultura, en donde los autores Sengupta y Lee, identifican el número de limones inmaduros bajo condiciones de temperatura ambiente. Los inputs que utilizan son *imágenes de la aspereza, brillo, suavidad, fineza del fruto, entre otros*. Son capaces de clasificar correctamente al 80.4% de los frutos (Sengupta y Lee, 2014). Para ver más aplicaciones en la agricultura, leer (Liakos y col., 2018).

Dado lo anterior, queremos señalar que estas metodologías pueden ser aplicables en diversas industrias y campos, los cuales pueden ayudar a aumentar la precisión de los problemas de *clasificación/predicción*.

Ahora, con respecto a literatura sobre ahorro y finanzas, el autor Lusardi utiliza los datos de la primera *wave* de la HRS (*Health and Retirement Study*), la cual es una encuesta que contiene información de individuos Estadounidenses nacidos en el año 1931-1941, en donde examina los hogares que se encuentran cercanos a la edad de jubilación. La HRS contiene una batería de preguntas con respecto a la salud futura, longevidad, precio de los hogares y algunas variables macroeconómicas. Lo encontrado por el autor radica en que gran parte de los hogares no han pensado sobre cómo planear la jubilación, de hecho, la forma en que los individuos aprenden a organizar su jubilación es mediante situaciones desagradables, como problemas financieros y shocks negativos en su salud; además, esta falta de programación induce a una poca tenencia de riqueza y en portafolios que son menos probables en contener activos de alto retorno, como las acciones. Asimismo, la responsabilidad del ahorro y contribuir a la pensión se ha cargado principalmente al trabajador, lo cual presenta un peligro al bienestar del individuo ya que este no realiza una planificación adecuada de su jubilación (Lusardi, 2003).

En lo referente a la educación financiera, Campbell encuentra que los hogares cometen muchos errores de inversión. El autor enfatiza en que los individuos no participan en mercados de activos riesgosos, no diversifican portafolios riesgosos y son incapaces de ejercer opciones de refinanciamiento de créditos hipotecarios en momentos de



tasas a la baja. También encuentra efectos heterogéneos en los hogares, ya que las familias más pobres y con menor educación son más propensos de cometer errores en comparación con los hogares de mayores ingresos y más educados (Campbell, 2006).

Existe diversa literatura que establece correlaciones entre la educación financiera con diferentes comportamientos financieros. Los autores (Hilgert y col., 2003), utilizando los datos de *University of Michigan's monthly Surveys of Consumers*, la cual mide los cambios en las actitudes y expectativas de los consumidores con respecto a decisiones financieras, presentan que existe una fuerte relación entre el conocimiento financiero y diversas prácticas financieras como pagar las cuentas a tiempo, hacer un seguimiento de los gastos, mantener un fondo de emergencia, inversiones diversificadas, entre otros. Van Rooij, Lusardi y Alessie examinan la relación entre el conocimiento financiero y la planificación de la jubilación utilizando datos de los países bajos *De Nederlandsche Bank (DNB) Household Survey*, en donde sus hallazgos radican en que el analfabetismo financiero está altamente generalizado en la población, además que la educación financiera tiene una fuerte correlación con la educación escolar y también manifiestan que la causalidad va desde la educación financiera hacia la preparación de la jubilación, y no al revés (van Rooij y col., 2011).

Con respecto a literatura de IA aplicada a finanzas, podemos encontrar tópicos tales como fraude financiero, la predicción de precios de activos, análisis en el comportamiento de inversores, optimización de portafolios, entre otros. Por ejemplo, En el trabajo de (Ravisankar y col., 2011), los autores utilizan datos de 202 compañías chinas y diversas metodologías como *Support Vector Machine*, *Redes neuronales*, entre otros, para predecir si encuentran fraude (malversación) en sus estados financieros. La predicción de fraude financieros es importante debido a que puede ahorrar una cantidad considerable de dinero ya que se evita el uso inapropiado de fondos.

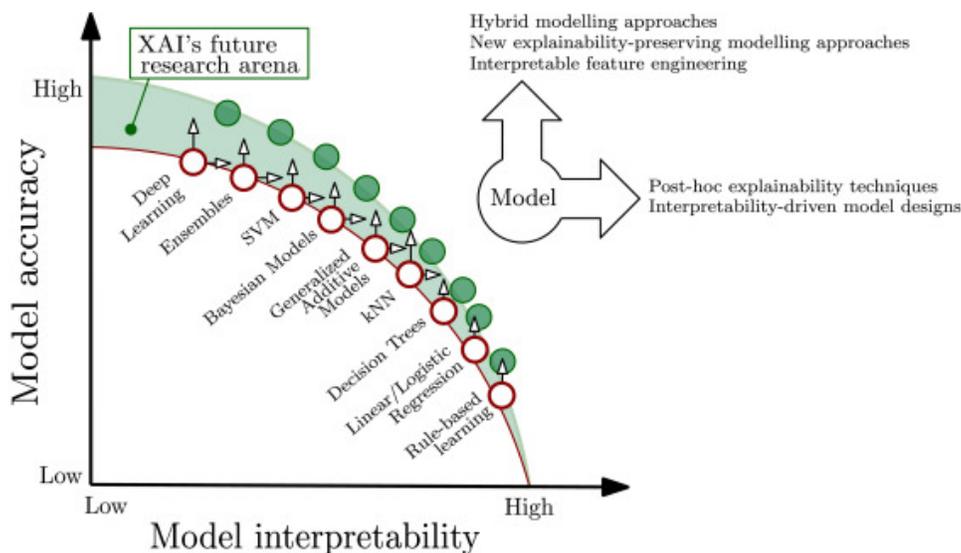
El artículo de (Nag y Mitra, 2002) tiene un enfoque en la predicción de *tipos de cambio* diario. Lo que encuentran los autores es que un tipo de modelo de *red neuronal* tiene un mejor performance en comparación con los modelos estadísticos clásicos (*arch*, *garch*). Para mayor detalle del uso del *Machine Learning* en finanzas, leer (Goodell y col., 2021).

Dentro del mundo del *Machine Learning*, existe un *Trade-Off* entre la capacidad predictiva de los algoritmos y la interpretabilidad de los resultados obtenidos. Por ejemplo, Las *Redes Neuronales* alcanzan una alta capacidad predictiva, pero la explicación de como llegaron a los resultados no es clara, a diferencia del algoritmo de *árboles de decisión*, en donde se puede observar el camino que sigue el algoritmo para poder llegar a la predicción final. Esto puede verse claramente en la figura 1, la cual



entrega diferentes modelos y la forma en como se comportan en interpretabilidad y *accuracy*.

Figura 1: Trade-off entre interpretabilidad y performance del algoritmo. Fuente: Barredo Arrieta y col., 2020.



Dado esto, el algoritmo final que se utiliza para poder predecir si el individuo sabe o no cuanto tiene acumulado en su cuenta individual de AFP corresponde a *Catboost*, el cual utiliza *Gradient Boosting* en los árboles de decisión (Dorogush y col., 2018). *Gradient Boosting*, en resumidas cuentas, utiliza los residuos de cada predicción, para ir agregándolos en las predicciones venideras con el objetivo de ir disminuyendo el error (Residuo) que obtienen las predicciones. Este proceso se repite muchas veces, por eso el nombre de Gradient Boosting.

Además, este algoritmo logra trabajar con variables categóricas de forma más armónica, a diferencia de otros algoritmos, en donde el pre-procesado para trabajar con variables categóricas debe ser más extenso.

Para la literatura sobre *Catboost*, existen diversas aplicaciones en donde se ha utilizado este algoritmo. Por ejemplo, (Hancock y Khoshgoftaar, 2020) utilizan *Catboost* para detectar fraudes en “Medicare” (Programa de Seguridad Social en Estados Unidos). Lo encontrado por estos autores es que este algoritmo supera a *XgBoost* en la detección de fraudes de “Medicare”, especialmente cuando se agregan variables relacionadas con la seguridad social.



Los autores (Jabeur y col., 2021), al comparar Catboost con otros 6 modelos, encuentran que este método tiene un mejor performance que los otros modelos para poder predecir una bancarrota corporativa. Por último, (Zhang y Fleyeh, 2019) han que Catboost es capaz de predecir de mejor forma los precios de electricidad de corto plazo en comparación con modelos de MLP (“Multy-Layer Perception”) o SVM (“Support Vector Machine”) en relación al performance de “Error Porcentual Absoluto Medio ” (MAPE).

Para poder entender las decisiones que toma nuestro algoritmos de *IA*, utilizamos el aporte que hace cada variable en la predicción de cada individuo. Para poder realizar esto, utilizamos *Shap Values* (Lundberg y Lee, 2017).

Shap utiliza un enfoque de teoría de juegos para poder explicar el output de modelos de *ML*. De forma más específica, *Shap* calcula la importancia de cada *feature* comparando lo que predice el modelo con y sin el *feature*. Esta comparación la hace realizando permutaciones, para que de esta manera las comparaciones sean justas y no dependan del orden en el cual fueron elegidas (Tseng, 2018).

El *Shap Value* de un *feature* para una observación se calcula de la siguiente forma (Mazzanti, 2020):

$$Shap_{feature}(x) = \sum_{set: feature \in set} \left[|set| \times \binom{F}{|set|} \right]^{-1} [Predict_{set}(x) - Predict_{set \setminus feature}(x)] \quad (1)$$

Donde x es una observación, F es el total de *features* y set es un subconjunto de F . La parte más a la derecha de la ecuación calcula como sería la predicción con y sin el *feature*, con lo cual realiza la diferencia. Esto sería el efecto del *feature* en la predicción. El valor elevado a la inversa corresponde a los pesos para poder realizar el promedio ponderado (se realizan permutaciones para poder evaluar todos los casos en donde se incluye o no el *feature*).

La escala del resultado *Shap Value* se encuentra en *Log-odds*, esto es simplemente otra forma de expresar una probabilidad. Matemáticamente,

$$Log\ odds = Ln \left(\frac{Pr(Si)}{1 - Pr(Si)} \right) \quad (2)$$



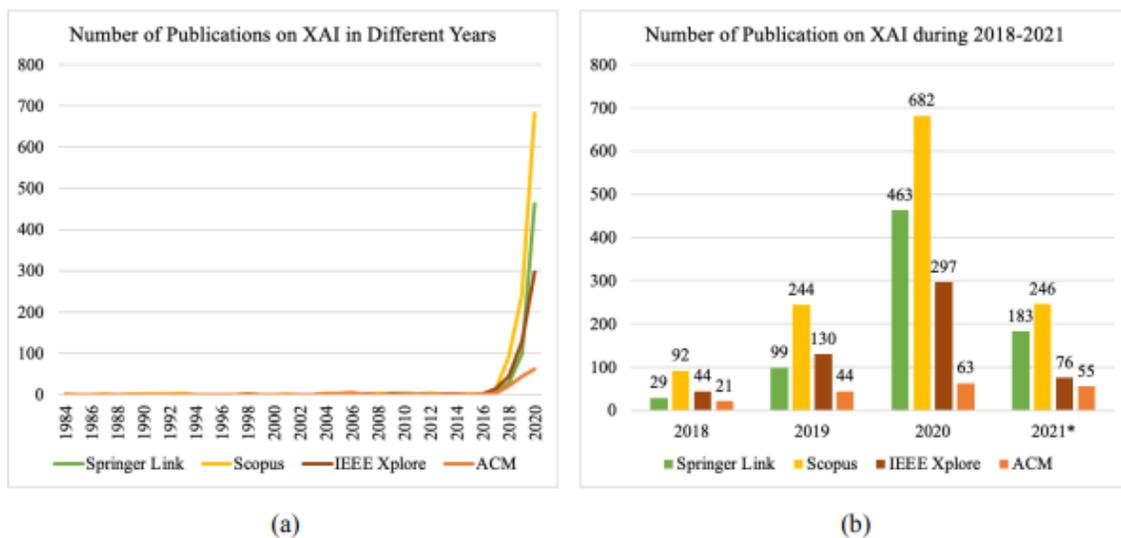
Luego, el valor en probabilidad es:

$$Pr = \frac{e^{\text{Log odds}}}{1 + e^{\text{Log odds}}} \quad (3)$$

Por lo tanto, cuando el valor de Log-odds es 0, la probabilidad predicha es de 50 %. Log-odds > 0 implica una probabilidad mayor al 50 %. Lo contrario ocurre cuando Log-odds es < 0, siendo una probabilidad menor al 50 %.

Con respecto al uso de *Explainable Artificial Intelligence (XAI)* en la literatura, esto se ha desarrollado principalmente en los últimos años. En la figura 2 se muestra el número de publicaciones relacionadas a XAI a lo largo del tiempo (panel a) y un zoom a los últimos años en donde explota la búsqueda de XAI en diferentes bases de datos bibliográficas (panel b). Se puede ver claramente que el auge en el número de publicaciones sobre XAI parte aproximadamente en el año 2016.

Figura 2: Literatura con Tópico XAI en el Tiempo. Fuente: Islam y col., 2022

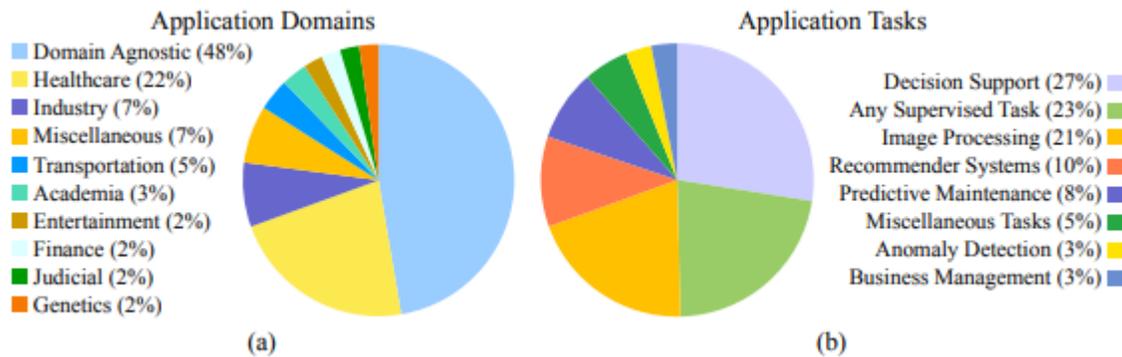


Diversos artículos han utilizado el concepto de *XAI*, en especial en el área de la medicina, en donde encontrar las causas de diversas enfermedades se ha vuelto un imperativo al momento de emplear herramientas de IA. Esto se ve reflejado en la figura 3, en donde el primer área de uso corresponde a "Domain Agnostic", lo cual



corresponde a un tópico explicado sin tener en cuenta ejemplos de ningún dominio en específico. Luego viene la industria de la salud. Con respecto a las tareas en específico, la principal es para el soporte de decisiones.

Figura 3: Literatura con Tópico XAI en el Tiempo. Fuente: Islam y col., 2022



3. Metodología, Pre-Proceso y Estadística Descriptiva

Los datos utilizados en el trabajo fueron obtenidos de la *Encuesta de Protección Social* (EPS) año 2015. Dentro de esta base, se encuentra nuestra variable que queremos predecir, la cual pregunta a los individuos *¿Sabe usted cuánto dinero hay acumulado en su cuenta individual?*, a lo cual tiene 2 respuestas posibles: *Si* y *No*.

Luego de limpiar la base de datos, contamos con un total de 9.671 observaciones. La distribución de nuestra variable *Target* se encuentra en la tabla 1. EL porcentaje de distribución de los individuos es aproximadamente de un 66% para los individuos que responden que **No** saben cuanto tienen acumulado en su cuenta de AFP y un 34% para los que **Sí** saben.



Tabla 1: Tabulación variable *¿Sabe usted cuánto dinero hay acumulado en su cuenta individual?*.

	Frecuencia	Porcentaje
No	6.338	65,54 %
Sí	3.333	34,46 %
Total	9.671	100 %

Las variables utilizadas para predecir son en total 36 variables, de las cuales 32 son variables categóricas y 4 son numéricas (números enteros o continuos).

Las variables categóricas se encuentran en la tabla 2 y las numéricas en la tabla 3. La estadística descriptiva de los *features* se encuentra en la tabla 7, en la sección de Anexos.



Tabla 2: Variables Categóricas Predictoras (Features Categóricos)

Variable
1 Sexo
2 ¿Pertenece o es descendiente de un pueblo indígena?
3 ¿Cuál es el nivel más alto alcanzado de educación?
4 Tipo de vivienda del entrevistado
5 ¿Tiene ahorros para la vivienda (banco)?
6 Usted o algún miembro de su hogar, ¿es dueño o socio de algún negocio o empresa?
7 Usted, ¿posee cuenta corriente?
8 Usted, ¿posee línea de crédito bancaria?
9 Usted, ¿posee préstamo de consumo bancario?
10 Usted, ¿posee préstamo de consumo en financieras?
11 Usted, ¿posee deudas educacionales?
12 ¿Usted tiene seguro de vida?
13 ¿Se encuentra cotizando actualmente?
14 ¿Sabe usted qué porcentaje de su ingreso imponible le descuentan?
15 En los últimos 12 meses, ¿ha recibido alguna cartola de su AFP?
16 ¿Conoce o ha escuchado de los multifondos?
17 ¿Sabe cuántos tipos de fondos existen?
18 Con respecto a su jubilación, ¿usted cree que dejará de trabajar?
29 ¿Conoce cuáles son las distintas modalidades de pensión por vejez?
20 ¿No ha recibido alguna proyección de su pensión?
21 ¿Conoce o ha escuchado hablar de la pensión básica solidaria?
22 ¿Conoce o ha escuchado hablar del aporte previsional solidario de vejez o APS?
23 ¿Ha escuchado hablar del seguro de cesantía?
24 ¿A qué sistema previsional de salud pertenece?
25 ¿Usted es cotizante o carga familiar?
26 En los últimos 2 años, ¿usted ha concurrido a un centro de salud?
27 En los últimos 12 meses, ¿ha solicitado licencias medicas?
28 ¿Tiene usted algún tipo de discapacidad?
39 Estado civil actual
30 ¿Cuál es la imagen que tiene de las AFP?
31 ¿Cuál es su imagen del sistema de pensiones de Chile?
32 En el caso de que existiese una AFP estatal, ¿usted está de acuerdo o en desacuerdo?

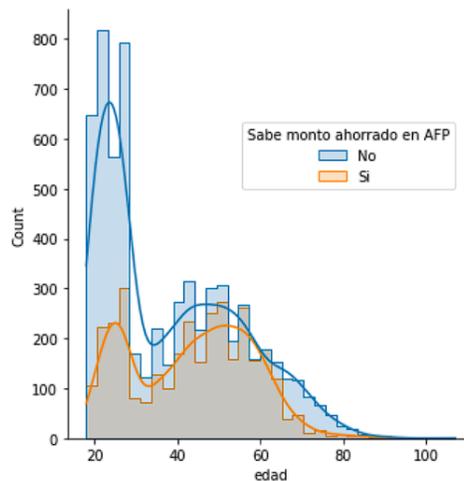


Tabla 3: Variables Numéricas Predictoras (Features Numéricos)

Variable
1 ¿En cuántas personas de su hogar, hijos o menores a carga debe gastar anualmente?
2 Edad
3 Ingreso mensual promedio de su hogar en los últimos 12 meses
4 En total, ¿cuántos hijos nacidos vivos o adoptados ha tenido usted?

Al observar la distribución de la edad de los individuos y diferenciando por si saben o no el monto ahorrado en su AFP, (figura 4) podemos ver que la densidad para individuos que responden *No* es mucho más concentrada en individuos más jóvenes. Por otro lado, para los individuos que responden *Si*, se puede ver que en la edad existe una distribución bimodal, concentrándose en las personas entre 20-30 y 50-60 años.

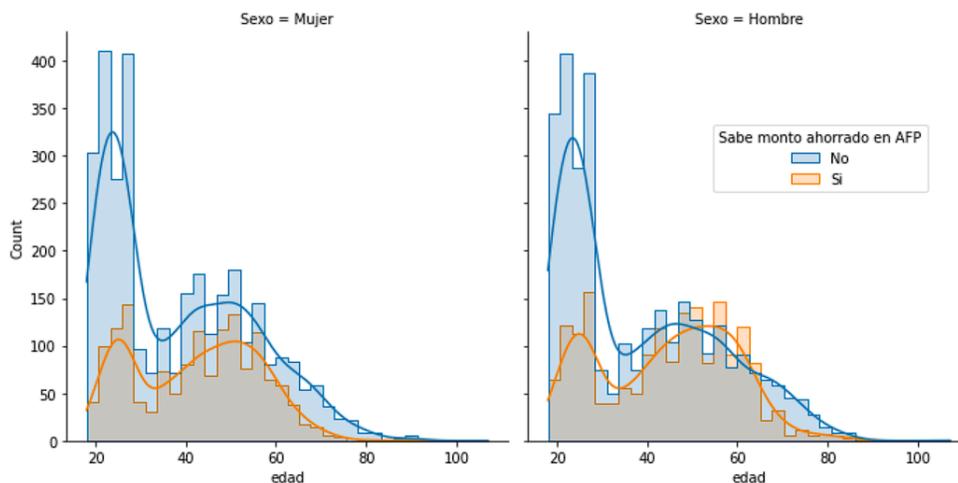
Figura 4: Distribución de la edad diferenciado por si saben el monto ahorrado en AFP.



Sin embargo, al desagregar por sexo del individuo (figura 5), si bien la forma de las distribuciones por lo general se mantienen, en el caso de los hombres la respuesta *Si* logra tener una densidad mayor en la edad cercana a los 60 años. Es decir, los hombres al acercarse a su edad de jubilación, es más probable que hayan observado cuanto dinero tienen en su cuenta de AFP.



Figura 5: Distribución de la edad diferenciado por si saben el monto ahorrado en AFP y género.



4. Resultados

En esta sección se entregan los resultados obtenidos con el algoritmo *Catboost*. Se entrega la matriz de confusión y el reporte de clasificación, el cual entrega el *accuracy*, *Precision* y *Recall* del modelo.

Además, tal como se ha mencionado anteriormente en el trabajo, los algoritmos de Machine Learning tienen un comportamiento de caja negra, sin embargo, la literatura ha hecho un avance en este aspecto y ha podido realizar diversas técnicas para poder abrir estas cajas negras y dar una mayor explicación a cómo los algoritmos de ML son capaces de entregar la predicción que realizan. Existen diversas librerías y/o algoritmos que abordan este punto, de los cuales utilizaremos *Shap Values* y *Shapash*. Con estos dos métodos podremos saber cuales fueron las variables más importantes dentro del modelo como un todo y también para cada individuo, de hecho, se puede conocer la manera en la cual los individuos tomaron sus decisiones. Esto quedará más claro en los gráficos de decisión en esta sección.

4.1. Estadísticas de Clasificación.

En primer lugar, es importante mencionar algunos conceptos para entender qué tan bueno es el algoritmo que se está utilizando. Dentro de estos, los más conocidos y/o importantes son:



- *Accuracy*: Es la proporción de las predicciones que fueron correctamente clasificadas. $Accuracy = \frac{N^{\circ} \text{ de clasificaciones correctas}}{\text{Total de Clasificaciones}}$
- *Precision*: Es la proporción de los resultados predichos como positivos que fueron correctamente clasificados. $Precision = \frac{\text{Verdadero Positivo}}{\text{Verdadero Positivo} + \text{Falso Positivo}}$
- *recall*: Es la proporción del total de positivos que fueron correctamente clasificados. $Recall = \frac{\text{Verdadero Positivo}}{\text{Verdadero Positivo} + \text{Falso Negativo}}$
- *F1-score*: Es un promedio ponderado entre *Precision* y *Recall*. $F_1 = \frac{2}{recall^{-1} + precision^{-1}}$

La tabla 4 muestra el reporte de clasificación para la muestra de Testing. Mediante Catboost se logra alcanzar un *Accuracy* del 74 % para la muestra de testing. La tabla 5 muestra los resultados de los otros modelos ordenados por *Accuracy*, de la cual se desprende que Catboost lo hace mejor en términos de *Accuracy*, seguido por *Random Forest* y *Logistic Regression*. Se puede observar que los valores cambian con respecto a las tablas de esta sección. Esto se debe a que la comparación con otros modelos se hizo utilizando *Cross-Validation* para obtener la estadística de las muestras de testeo y para el reporte de clasificación individual de catboost se utiliza una división de la muestra en 80-20 (*split training-test*).

De la tabla 5 podemos notar que *Catboost* tiene un mucho mejor desempeño en *Precision* en comparación con *Logistic Regression*, es decir, Catboost genera menos *falsos positivos*, por lo cual tiene un mejor performance para las clasificaciones observadas como positivas. Caso contrario ocurre con el *Recall*, en donde *Logistic Regression* genera menos *falsos negativos*. No obstante, tal como hemos señalado, Catboost se desempeña mejor en términos de *Accuracy* y prácticamente iguales en *AUC*.

Tabla 4: Reporte de Clasificación Muestras de Training y Testing.

Target	Training				Testing			
	Precision	Recall	F1-Score	Accurracy	Precision	Recall	F1-Score	Accurracy
No	0.82	0.79	0.80	0.81	0.86	0.73	0.79	0.74
Si	0.79	0.83	0.81		0.60	0.77	0.67	



Tabla 5: Comparación de estadísticas de clasificación con otros algoritmos.

Model	Accuracy	AUC	Recall	Prec.
Catboost	0.7616	0.8193	0.6221	0.6627
Random Forest Classifier	0.7572	0.8121	0.5958	0.6632
Logistic Regression	0.7524	0.8192	0.7449	0.6151
Decision Tree Classifier	0.6759	0.6464	0.5524	0.5272
Support Vector Machine	0.7375	0.0000	0.7544	0.5943
K Neighbors Classifier	0.6390	0.7144	0.7250	0.4831

La tabla 6 entrega el Cross Validation solo del modelo Catboost. De esta manera entregamos mayor robustez a nuestros resultados, ya que realizamos un training y testing en 10 muestras distintas. Se puede notar que los valores promedios son iguales a los de la tabla 5.

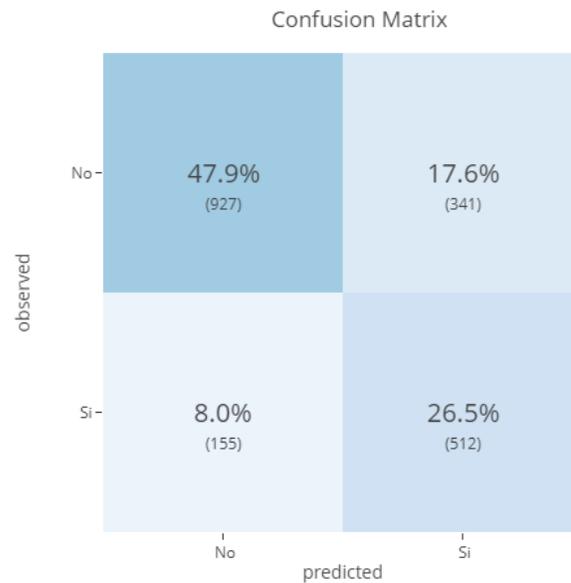
Tabla 6: Cross Validation Catboost model.

Fold	Accuracy	AUC	Recall	Prec.	F1
0	0.7434	0.8090	0.5882	0.6373	0.6118
1	0.7496	0.8229	0.6335	0.6364	0.6349
2	0.7543	0.8278	0.6109	0.6522	0.6308
3	0.7838	0.8309	0.7104	0.6767	0.6932
4	0.7527	0.8056	0.5973	0.6535	0.6241
5	0.7527	0.8058	0.5747	0.6615	0.6150
6	0.7745	0.8303	0.6244	0.6900	0.6556
7	0.7760	0.8263	0.6409	0.6845	0.6620
8	0.7698	0.8305	0.6227	0.6782	0.6493
9	0.7589	0.8038	0.6182	0.6570	0.6370
Mean	0.7616	0.8193	0.6221	0.6627	0.6414
SD	0.0128	0.0111	0.0353	0.0180	0.0233

La figura 6 entrega la matriz de Confusión para la muestra de Testing. Se puede observar que la mayor cantidad de observaciones se encuentra en los *Verdaderos Negativos*, con una cantidad de 927 observaciones. El total de predicciones para la muestra de Testing fue de 1935.



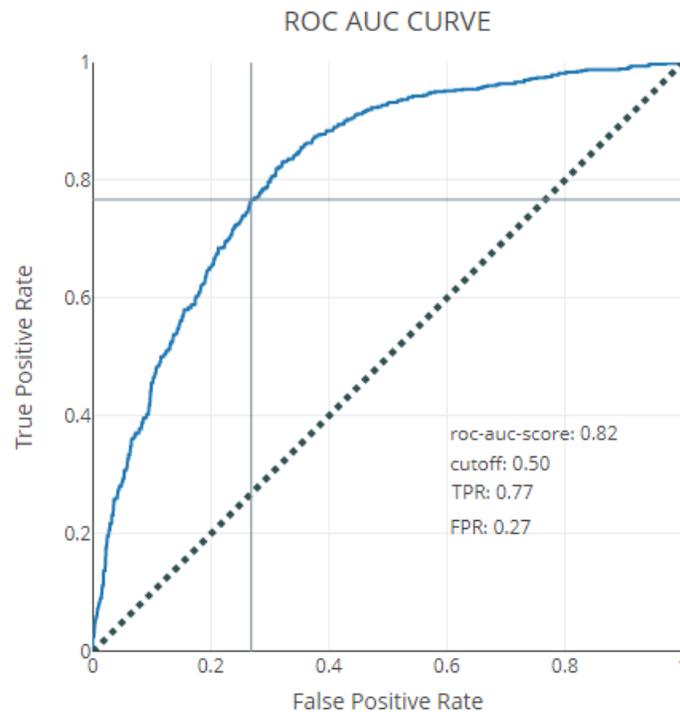
Figura 6: Matriz de Confusión



La figura 7 enseña la curva ROC-AUC (ROC: Receiver Operator Characteristic. AUC: Area Under the Curve.) del modelo. Con esta figura podemos ver un resumen del Trade-off entre el ratio de Falsos Positivos con el ratio Verdaderos Positivos cuando se cambia el umbral de clasificación (en el modelo se utiliza un umbral de $> 50\%$ para clasificar a un individuo como Positivo). La Línea punteada es donde el ratio de los Falsos Positivos es igual al Ratio de los Verdaderos Positivos. EL AUC obtiene un valor de 0.82 (entre más cercano a 1, mejor es el modelo).



Figura 7: Curva ROC-AUC



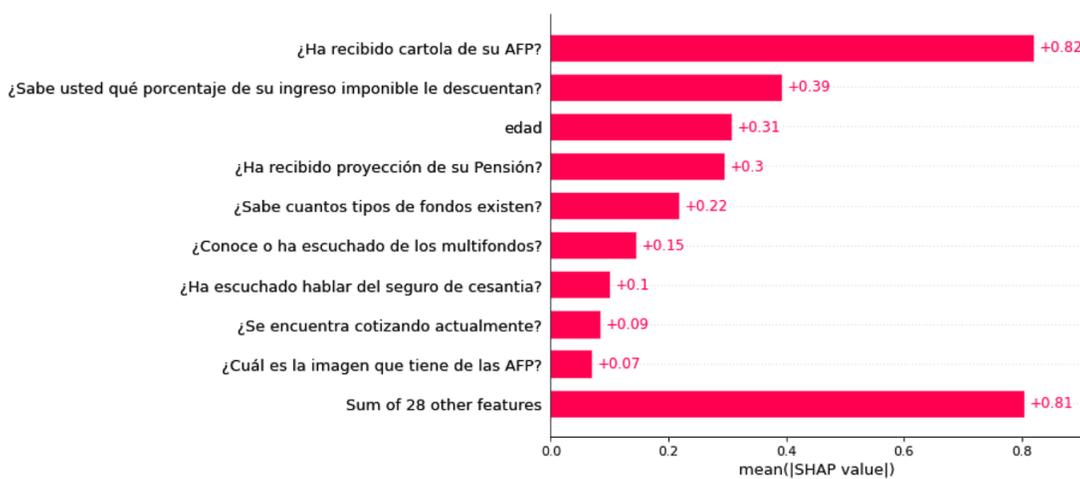
4.2. Análisis del modelo.

En esta sección se realiza un análisis macro del modelo, es decir, cuales son las variables más importantes para realizar la predicción y cómo se comportan estas variables.

Dada la introducción de *Shap* que fue vista en la sección de revisión de literatura, podemos calcular los *Shap Values* para todos los individuos y *features*. En La figura 8 se ordenan las variables por su impacto en la predicción (se considera el promedio del valor absoluto del *Shap Value*). Se puede desprender que para la muestra de Testing, las 3 variables que más contribuyen en promedio a la predicción son: *En los últimos 12 meses, ¿ha recibido alguna cartola de su AFP?*, *¿Sabe usted qué porcentaje de su ingreso imponible le descuentan?* y la *edad* del individuo.

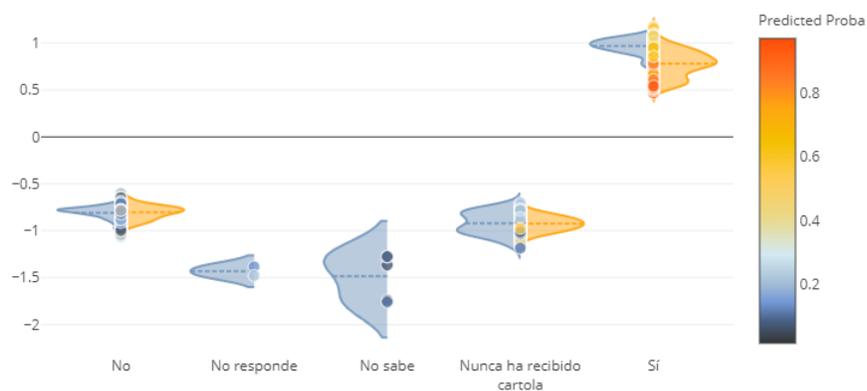


Figura 8: Feature Importance por Shap Values. Muestra de Testing.



Sin embargo, la figura 8 solo nos muestra la importancia de las variables para las predicciones, pero no la dirección de esta magnitud (impacto positivo o negativo en la predicción). Para observar la dirección, utilizamos los gráficos *Feature Contribution* de la librería *Shapash* (MAIF, 2020). Con este tipo de gráfico podemos observar el impacto que tienen los distintos valores de las variables en la predicción (estos gráficos también están en escala \log_odds). Los gráficos 9 y 10 entregan los *Feature Contribution* de las 2 variables más importantes según la figura 8.

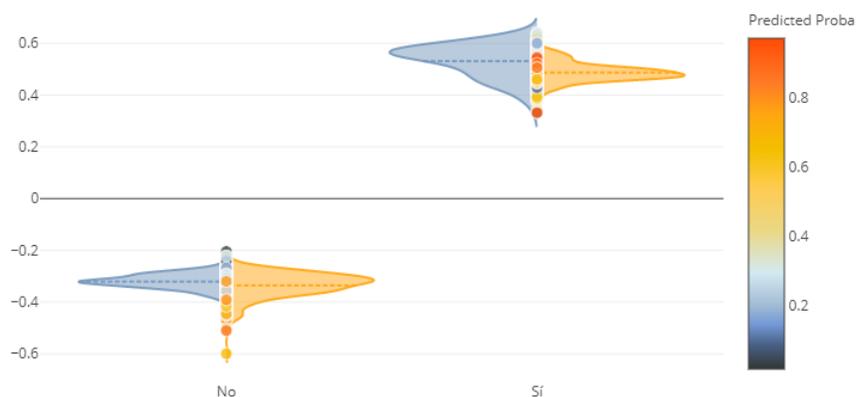
Figura 9: Feature Contribution variable: *En los últimos 12 meses, ¿ha recibido alguna cartola de su AFP?*





De la figura 9 se puede inferir que los individuos que responden que “Sí” a la pregunta, “¿ha recibido alguna cartola de su AFP?”, la predicción del modelo se inclina más por la opción de que el individuo **Sí** sepa cuánto tiene acumulado en su cuenta individual. Caso contrario ocurre con las demás respuestas, en donde se contribuye más por la opción de que el individuo **No** sepa cuanto hay acumulado en su cuenta individual.

Figura 10: Feature Contribution variable: ¿Sabe usted qué porcentaje de su ingreso imponible le descuentan?



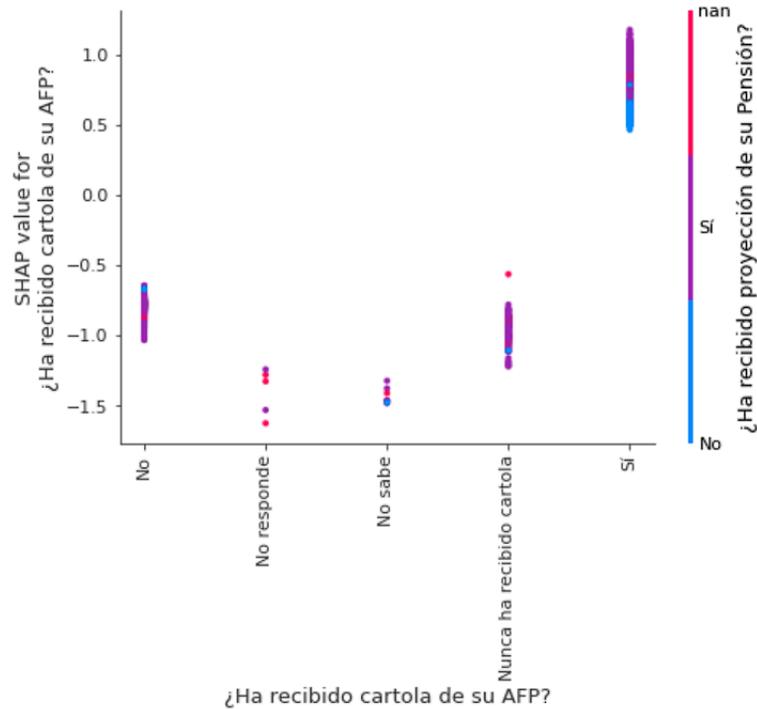
En la figura 10 se puede ver que los individuos que No saben cuanto porcentaje de su ingreso imponible le descuentan, el modelo se inclina a que estos individuos No sepan cuanto tienen acumulado en su cuenta individual de AFP. A diferencia de los individuos que sí saben cuanto es el porcentaje que les descuentan, en donde el efecto es el contrario.

Dado que existe una dispersión vertical en los gráficos de *Feature Contribution*, se puede inferir que existen efectos de interacción con otras variables. Para observar esta interacción, se realizan los gráficos *Dependence plot*, en donde otra variable es elegida para colorear y resaltar las posibles interacciones.

Por ejemplo, la figura 11 muestra que los individuos que responden que sí han recibido una cartola de la AFP, el efecto positivo sobre que sepan cuanto hay acumulado en su cuenta individual es mayor, en general, para individuos que sí han recibido proyección de su pensión.



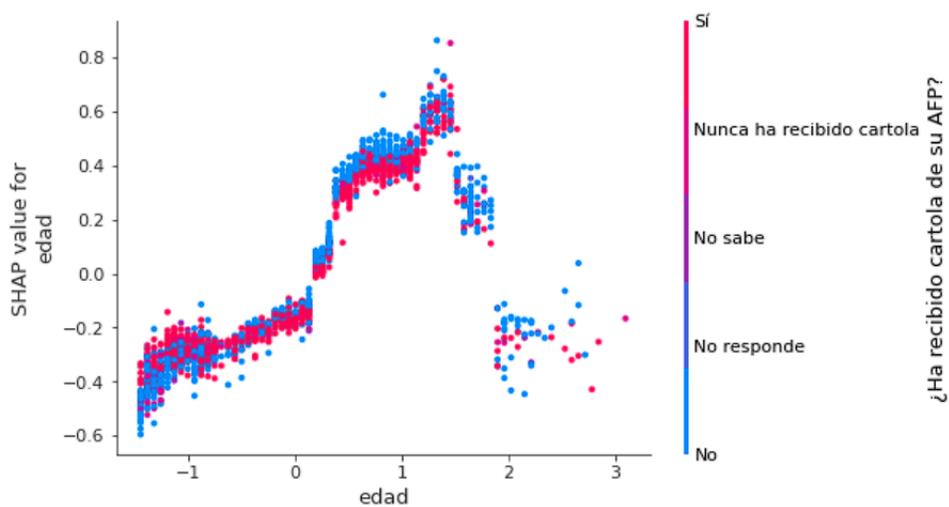
Figura 11: Dependence Plot variables ¿Ha recibido cartola de su AFP? y ¿Ha recibido proyección de su pensión?



La figura 12 muestra que al aumentar la edad, aumenta la probabilidad de que el output sea **Sí**, sin embargo, el efecto no es lineal, sino cóncavo ya que a partir edades más tardías, la contribución empieza a declinar. Además, se puede ver que los individuos de mayor edad en su mayoría han respondido que *no* han recibido cartola de su AFP en los últimos 12 meses.



Figura 12: Dependence Plot variables *edad* y *¿Ha recibido cartola de su AFP?*



Cabe mencionar que el gráfico 12, está con la edad estandarizada, y no se centra en el cero ya que la estandarización proviene de la muestra de Entrenamiento y el gráfico corresponde a la muestra de prueba. El promedio de edad en la muestra de Entrenamiento fue de 40,89 y la desviación estándar fue de 15,97.

4.3. Análisis de individuos.

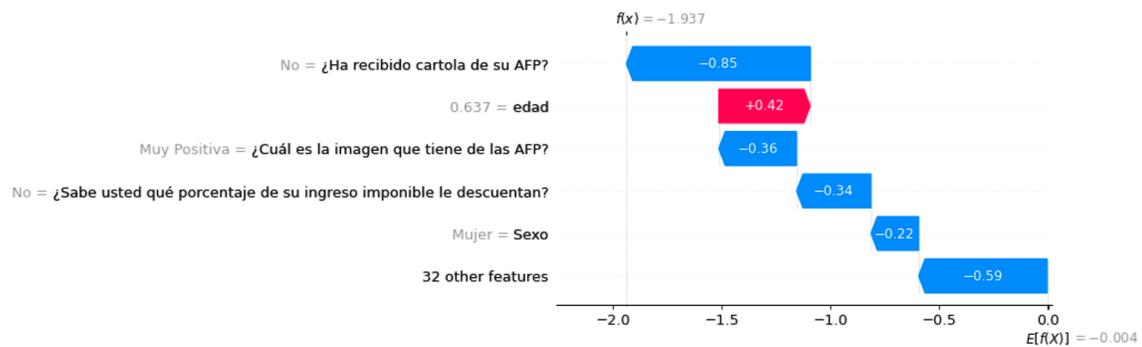
Luego del análisis del modelo, en esta sección haremos un enfoque a la manera en la cual algunos individuos obtuvieron su predicción. Se realiza este enfoque debido a que se quiere caracterizar a individuos que consideramos importantes en la muestra, tales como personas con cierta opinión de las AFP, su sexo, edad, entre otros.

En primer lugar, nos parece interesante ver el output que tienen los individuos que han sido más influenciados por la variable *¿Cuál es la imagen que tiene de las AFP?*. En la figura 13 se muestra la manera en la que se llegó a la predicción para el individuo en donde la *imagen que tiene de la AFP* obtiene el mayor valor absoluto en *shap value*. Este individuo tiene una imagen “*muy positiva*” de las AFP, sin embargo esto contribuyó (para este individuo) a una menor probabilidad de que el individuo **Sí** haya pensado en cómo financiar su vejez, lo cual es contrario a la intuición, ya que uno esperaría que individuos con una imagen positiva de las AFP, tendrían un impacto positivo en la probabilidad de que piensen en el financiamiento de la vejez (también es válida la intuición de que al tener una imagen positiva de las AFP, estarían menos



preocupados del monto ahorrado ya que este se encuentra delegado en una institución en la cual se tiene una buena imagen). Cabe mencionar que aunque este individuo fue el que tuvo el mayor *Shap Value* en valor absoluto para la variable *imagen que tiene de la AFP*, esta no fue la variable que más influyó en su predicción, siendo esta última la variable *¿Ha recibido cartola de su AFP?* en donde la respuesta de este individuo en particular es *No*.

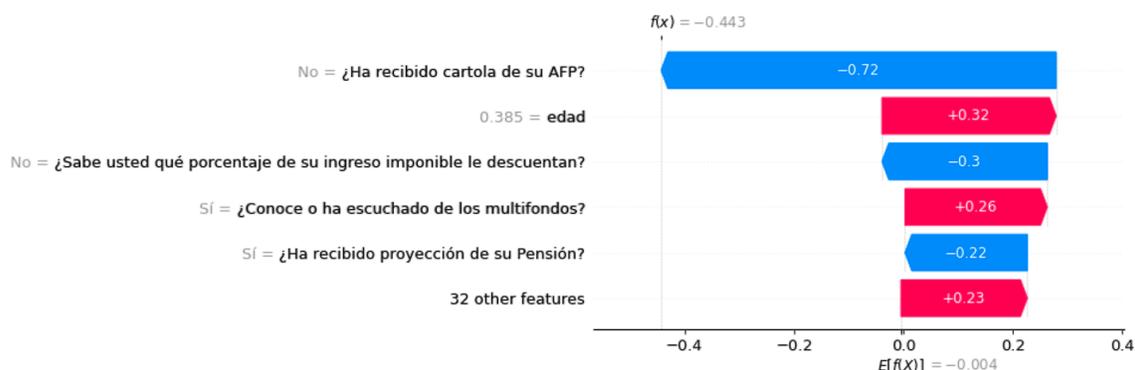
Figura 13: Individuo más influenciado por variable: *imagen que tiene de las AFP*



Luego, en la figura 14, se muestra a un individuo de sexo femenino, 47 años de edad, con Educación Media. Se muestra este individuo ya que respondió que *seguirá trabajando hasta que la salud se lo permita* en la variable “*Con respecto a su jubilación, ¿usted cree que dejará de trabajar?*”. Este individuo obtiene una predicción de **No** para el *target*. Se puede observar que la variable que más influye en la predicción es *¿Ha recibido cartola de su AFP?*, en donde la respuesta de este individuo para esta variable es “*No*”.



Figura 14: Individuo femenino, 20 años, Educación Media.



Tal como fue presentado en esta sección, podemos caracterizar a todos los individuos de la muestra, observando cuales fueron las variables más importantes para la predicción que otorgó el modelo. De esta forma, podemos “abrir” la caja negra que son los algoritmos de *ML* y tratar de explicar las decisiones que toman estos algoritmos.

4.4. Relevancia Económica

Luego del análisis global y local del modelo, en esta sección se procede a mostrar la forma en la cual individuos que **No** saben cuanto hay acumulado en su cuenta individual, podrían cambiar su respuesta a **Sí**. Esto con el fin de motivar políticas públicas que aumenten la probabilidad de que los individuos se preocupen del financiamiento de su vejez.

El procedimiento es el siguiente, en primer lugar se observan las variables que son susceptibles a un cambio en la respuesta que da el individuo. Por ejemplo, la edad del individuo no es posible cambiarla, pero sí es factible entregar al individuo la información de cuanto de su ingreso imponible le descuentan, variable que contamos en nuestra base de datos. En consecuencia, Las variables relevantes para hacer la política son primordialmente en las que los individuos tienen una falta de conocimiento, lo cual puede ser corregido mediante una entrega novedosa y de bajo costo de la información. Esto se podría reflejar, en la cartola de sueldo de los trabajadores (la cual es más probable que vean en comparación con la cartola de la AFP), en donde podría haber una sección con información del ahorro que ha hecho el individuo en el sistema de pensiones (monto acumulado, proyección de pensión, entre otros). Asimismo, esta misma cartola puede contener una sección que explique de



forma clara y concisa conceptos relacionados con el sistema de pensiones, tal es el caso de explicar que son los multifondos y cómo se componen, ahorro previsional voluntario (APV), modalidades de pensión, seguro de cesantía, entre otros. De esta manera se contribuiría a que el individuo tenga un mayor conocimiento del sistema y de su pensión.

Además, como vimos en los resultados del impacto de las variables, sabemos que los individuos que no saben el porcentaje que le descuentan del sueldo imponible influye a que tengan una menor probabilidad de saber cuanto es el monto que tienen ahorrado en su AFP. Una manera de reducir esta brecha sería nuevamente utilizando las cartolas de sueldos, en donde se podrían mostrar no solo los niveles de descuento en un monto monetario, sino también en porcentaje con respecto al total bruto. Por ejemplo, que en la cartola exista una columna en la cual aparezca el porcentaje del total bruto que se destina a cotización de AFP (11 % aprox.), el porcentaje que se destina a salud (7 % aprox.), impuestos, seguros, entre otros, para luego mostrar el total que se descontó en porcentaje (por ejemplo un individuo tuvo un descuento total del 24 % del sueldo imponible). De esta forma, el individuo sabría el total de descuentos que se le está realizando y la composición de este, contribuyendo a que el individuo tenga una mayor información de los descuentos a su sueldo imponible.

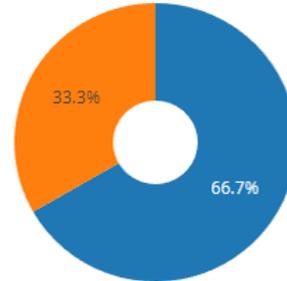
Para ilustrar lo discutido en esta sección, la figura 15 muestra un individuo el cual **No** sabe *cuanto porcentaje de su ingreso imponible le descuentan* y además tiene una probabilidad predicha de 33,3 %. Sin embargo, si se le entregara esta información, la probabilidad predicha pasa a ser del 56 %, lo cual es mayor al umbral y se clasificaría como una persona que sí sabe cuanto hay acumulado en su cuenta individual. Este es un ejemplo de cómo una política pública podría cambiar el output del individuo.



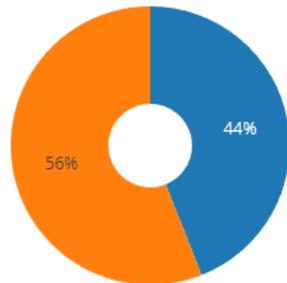
Figura 15: Antes y después de saber cuanto porcentaje de su ingreso imponible le descuentan. Individuo 14.634.

Prediction

label	probability	logodds
No	66.7 %	0.692
Si	33.3 %	-0.692



label	probability	logodds
No	44.0 %	-0.242
Si	56.0 %	0.242



4.5. Discusión.

En esta sección discutiremos nuestros resultados encontrados, entregando una mayor explicación de lo que descubrimos y dando posibles cambios que se deberían realizar en el futuro.

Nuestro principal foco en el trabajo es tratar de predecir *si un individuo sabe cuánto dinero hay acumulado en su cuenta individual de AFP*. Luego, intentamos averiguar cuáles fueron las variables más importantes para la respuesta que entregan los individuos y cómo se comportarían si cambiáramos la respuesta de alguna variable.

Dado nuestros resultados encontrados, obtenemos un buen desempeño en la predicción, sin embargo, el *“Recall”* puede ser bajo, pero no es un número con un desempeño pésimo. Esto se podría mejorar, tal como lo hicimos, con un rebalanceo de la data u



obteniendo una muestra mucho más grande.

Luego, una de las principales riquezas de nuestro trabajo, es que podemos mostrar cuáles son las variables más importantes en nuestro modelo (tal como hemos discutido a lo largo de este trabajo), e incluso observar cómo se comportan las variables para cada individuo. Esto es un tópico relativamente nuevo en el mundo de la Inteligencia Artificial. Si bien, se podría argumentar que ya existen algoritmos “tradicionales” como *Regresión Logística* que ya pueden hacer esto, estos algoritmos utilizan una mayor cantidad de tiempo para lograr los resultados y, lo más importante, tienen una menor capacidad de *predicción*, lo cual es uno de los componentes que más nos interesan.

Por último, debemos mencionar que este tipo de trabajos, es decir, tratar de descubrir qué es lo que hacen o cómo toman las decisiones este tipo de algoritmos de caja negra, han tenido en el último tiempo una mayor cantidad de desarrollo, por lo cual en los siguientes años se pueden emitir nuevos artículos que explayan aún mejor los algoritmos de Inteligencia Artificial.

5. Conclusión

El objetivo de este trabajo ha sido tratar de predecir si un individuo sabe cuánto hay acumulado en su cuenta individual de AFP. Dado los resultados encontrados, nuestro trabajo pudo predecir correctamente el 74% de las predicciones realizadas. Cabe recordar que nuestro enfoque es la predicción y no la causalidad.

Las variables más importantes que hemos encontrado para realizar las predicciones han sido “¿en los últimos 12 meses, ha recibido alguna cartola de su AFP?”, “¿sabe usted que porcentaje de su ingreso imponible le descuentan?” y “edad”.

Dado este trabajo, se recomienda como política pública ir especialmente a los individuos que responden (o tienen un valor predicho) que **No** saben cuanto monto hay acumulado en su cuenta individual de AFP, en especial en individuos más jóvenes, en donde los cambios en los comportamientos de ahorro pueden entregar un gran impacto en el largo plazo, mejorando sustancialmente su pensión. Tal como se discutió en la sección de relevancia económica, una forma de política pública es que las cartolas no estén separadas, sino que se entreguen en conjunto, de esta forma se hace más probable que los individuos puedan observar información relacionada a su ahorro en el sistema de pensiones. Además, se puede agregar a la cartola información del sistema de pensión, tal como las modalidades de jubilación, seguro de cesantía,



entre otros. Por otro lado, en la misma cartola se podrían entregar los porcentajes correspondientes a cada descuento (AFP, salud, etc.), con lo cual los individuos sabrían el monto total de descuento y su composición, relativo al sueldo bruto.

Para investigación futura, es recomendable volver a revisar los algoritmos que tratan de explicar las decisiones que realizan los programas de Machine Learning, esto debido a que es un campo relativamente nuevo, con lo cual se pueden emitir próximas herramientas que nos ayuden aún más en el trabajo de abrir la caja negra.



Referencias

- Barredo Arrieta, A., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbadó, A., Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R. & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82-115. <https://doi.org/https://doi.org/10.1016/j.inffus.2019.12.012>
- Behrman, J. R., Mitchell, O. S., Soo, C. K. & Bravo, D. (2012). How financial literacy affects household wealth accumulation. *American Economic Review*, 102(3), 300-304.
- Campbell, J. Y. (2006). Household finance. *The journal of finance*, 61(4), 1553-1604.
- Clark, R. L., Morrill, M. S. & Allen, S. G. (2012). THE ROLE OF FINANCIAL LITERACY IN DETERMINING RETIREMENT PLANS. *Economic Inquiry*, 50(4), 851-866. <https://doi.org/https://doi.org/10.1111/j.1465-7295.2011.00390.x>
- Dorogush, A. V., Ershov, V. & Gulin, A. (2018). CatBoost: gradient boosting with categorical features support. *CoRR*, abs/1810.11363. <http://arxiv.org/abs/1810.11363>
- Goodell, J. W., Kumar, S., Lim, W. M. & Pattaik, D. (2021). Artificial intelligence and machine learning in finance: Identifying foundations, themes, and research clusters from bibliometric analysis. *Journal of Behavioral and Experimental Finance*, 32, 100577.
- Hancock, J. & Khoshgoftaar, T. M. (2020). Medicare fraud detection using catboost. *2020 IEEE 21st international conference on information reuse and integration for data science (IRI)*, 97-103.
- Hilgert, M. A., Hogarth, J. M. & Beverly, S. G. (2003). Household financial management: The connection between knowledge and behavior. *Fed. Res. Bull.*, 89, 309.
- IBM. (2020). *What is machine learning?* <https://www.ibm.com/cloud/learn/machine-learning#toc-what-is-ma-qhM6PX35>
- Islam, M. R., Ahmed, M. U., Barua, S. & Begum, S. (2022). A systematic review of explainable artificial intelligence in terms of different application domains and tasks. *Applied Sciences*, 12(3), 1353.
- Jabeur, S. B., Gharib, C., Mefteh-Wali, S. & Arfi, W. B. (2021). CatBoost model and artificial intelligence techniques for corporate failure prediction. *Technological Forecasting and Social Change*, 166, 120658. <https://doi.org/https://doi.org/10.1016/j.techfore.2021.120658>



- Klapper, L. F., Lusardi, A. & Panos, G. A. (2012). *Financial Literacy and the Financial Crisis* (Working Paper N.º 17930). National Bureau of Economic Research. <https://doi.org/10.3386/w17930>
- Liakos, K. G., Busato, P., Moshou, D., Pearson, S. & Bochtis, D. (2018). Machine Learning in Agriculture: A Review. *Sensors*, 18(8). <https://doi.org/10.3390/s18082674>
- Lundberg, S. M. & Lee, S.-I. (2017). A Unified Approach to Interpreting Model Predictions. En I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan & R. Garnett (Eds.), *Advances in Neural Information Processing Systems 30* (pp. 4765-4774). Curran Associates, Inc. <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
- Lusardi, A. (2003). *Planning and saving for retirement* (inf. téc.). Working paper. Dartmouth College.
- MAIF. (2020). *Shapash*. <https://github.com/MAIF/shapash>
- Mazzanti, S. (2020). *SHAP Values Explained Exactly How You Wished Someone Explained to You*. <https://towardsdatascience.com/shap-explained-the-way-i-wish-someone-explained-it-to-me-ab81cc69ef30>
- McCarthy, J. (2004). WHAT IS ARTIFICIAL INTELLIGENCE? *Standford*.
- Nag, A. K. & Mitra, A. (2002). Forecasting daily foreign exchange rates using genetically optimized neural networks. *Journal of Forecasting*, 21(7), 501-511.
- Ravisankar, P., Ravi, V., Rao, G. R. & Bose, I. (2011). Detection of financial statement fraud and feature selection using data mining techniques. *Decision support systems*, 50(2), 491-500.
- Sengupta, S. & Lee, W. S. (2014). Identification and determination of the number of immature green citrus fruit in a canopy under different ambient light conditions [Image Analysis in Agriculture]. *Biosystems Engineering*, 117, 51-61. <https://doi.org/https://doi.org/10.1016/j.biosystemseng.2013.07.007>
- Shailaja, K., Seetharamulu, B. & Jabbar, M. (2018). Machine learning in healthcare: A review. *2018 Second international conference on electronics, communication and aerospace technology (ICECA)*, 910-914.
- Tseng, G. (2018). *Interpreting complex models with SHAP values*. <https://medium.com/@gabrielteng/interpreting-complex-models-with-shap-values-1c187db6ec83>
- Turing, A. M. (2009). Computing machinery and intelligence. *Parsing the turing test* (pp. 23-65). Springer.
- van Rooij, M. C., Lusardi, A. & Alessie, R. J. (2011). Financial literacy and retirement planning in the Netherlands [Financial Capability]. *Journal of Economic Psychology*, 32(4), 593-608. <https://doi.org/https://doi.org/10.1016/j.joep.2011.02.004>



- Venkatesan, E. & Velmurugan, T. (2015). Performance analysis of decision tree algorithms for breast cancer classification. *Indian Journal of Science and Technology*, 8(29), 1-8.
- Zhang, F. & Fleyeh, H. (2019). Short term electricity spot price forecasting using catboost and bidirectional long short term memory neural network. *2019 16th International Conference on the European Energy Market (EEM)*, 1-6.



6. Anexos

Tabla 7: Estadística descriptiva Features.

Variable	Obs	Mean	Std. Dev.
Sexo	9,671	1.4924	.4999681
¿Pertenece o es descendiente de un pueblo indígena?	9,671	10.23307	8.098455
¿cuál es el nivel más alto alcanzado de educación?	9,671	7.151174	4.744441
Edad	9,671	40.0304	16.03946
Ingreso mensual promedio de su hogar en los últimos 12 meses	9,671	564191.8	785645.5
Tipo de vivienda del entrevistado	9,671	1.432634	1.22571
¿Tiene ahorros para la vivienda (banco)?	9,671	1.92524	.5158249
Usted o algún miembro de su hogar, ¿es dueño o socio de algún negocio o empresa?	9,671	2.931445	.4427455
¿Posee cuenta corriente?	9,671	1.903112	.5082786
¿Posee línea de crédito bancaria?	9,671	1.952332	.5531731
¿posee préstamo de consumo bancario?	9,671	1.937752	.5290454
¿posee préstamo de consumo en financieras?	9,671	2.006514	.4395859
¿posee deudas educacionales?	9,671	1.968979	.6867556
¿En cuántas personas de su hogar, hijos o menores a carga debe gastar anualmente?	9,671	.7462517	1.318208
¿Usted tiene seguro de vida?	9,671	1.950367	.7728709
¿Se encuentra cotizando actualmente?	9,671	1.465929	.4988636
¿Sabe usted qué porcentaje de su ingreso imponible le descuentan?	9,671	1.6626	.4774179
En los últimos 12 meses, ¿ha recibido alguna cartola de su AFP?	9,671	1.666322	.8076816
¿Sabe usted cuanto hay acumulado en su cuenta individual?	9,671	1.655361	.4752749
¿Conoce o ha escuchado de los multifondos?	9,671	1.652363	.4762447
¿Sabe cuántos tipos de fondos existen?	9,671	2.157274	1.786013
Con respecto a su jubilación, ¿usted cree que dejará de trabajar?	8,624	7.058558	18.60344
¿Conoce cuáles son las distintas modalidades de pensión por vejez?	9,671	2.128322	1.239381
¿Ha recibido alguna proyección de su pensión?	8,631	1.140887	.3479256
¿Conoce o ha escuchado hablar de la pensión básica solidaria?	9,671	1.77841	1.244725
¿Conoce o ha escuchado hablar del aporte previsional solidario de vejez o APS?	9,671	2.142695	1.286598
¿Ha escuchado hablar del seguro de cesantía?	9,671	1.335643	.8422961
¿A qué sistema previsional de salud pertenece?	9,671	4.426636	10.38422
¿Usted es cotizante o carga familiar?	9,343	1.954618	2.404163
En los últimos 2 años, ¿usted ha concurrido a un centro de salud?	9,671	1.304622	.4649653
En los últimos 12 meses, ¿ha solicitado licencias medicas?	9,671	1.890187	.3126725
¿Tiene usted algún tipo de discapacidad?	9,671	1.982422	.42187
Estado civil actual	9,671	4.228105	3.591597
En total, ¿cuántos hijos nacidos vivos o adoptados ha tenido usted?	9,671	1.625685	1.765666
¿Cuál es la imagen que tiene de las AFP?	9,671	3.677283	.9685755
¿Cuál es su imagen del sistema de pensiones de Chile?	9,671	3.917382	.9248376
En el caso de que existiese una AFP estatal, ¿usted está de acuerdo o en desacuerdo?	9,671	2.09058	1.977711