



Universidad de Buenos Aires
Facultad de Ciencias Económicas
Escuela de Estudios de Posgrado



Carrera de Especialización en Métodos Cuantitativos para la gestión y análisis de datos en Organizaciones

Trabajo Final de Especialización

Un modelo de predicción de próxima compra aplicado a bancos

Aplicación de técnicas de minería de datos en un entorno de Python para mejorar la estrategia de marketing directo.

Autor: Kramarz, Brian Ariel

Diciembre 2019

Introducción

Durante las últimas décadas, los bancos se han convertido en uno de los tantos tipos de organizaciones que han comenzado a recopilar datos de sus clientes de forma sistemática. Si bien existen ciertos datos que podemos esperar que toda empresa tenga respecto de sus clientes, como por pueden ser ciertos datos personales como domicilio, datos de contacto, nacionalidad, entre otros, debido a su tipo de actividad, los banco tienen también mucha otra información más sensible como podría ser el estado civil, si un cliente tiene gente a cargo, e incluso sus ingresos periódicos. Es más, con el avance de los pagos por medios electrónicos, ya sea tarjeta de crédito y/o debido, pago mis cuentas y las nuevas billeteras virtuales, los bancos puede conocer que consumimos, cuando y en donde realizamos nuestras compras, y muchos más respecto de nuestros consumos sean estos habituales o no. Cabría preguntarse cuál, si alguna de toda esa información, los bancos podrían utilizarse para sus propios intereses, lo que sin dudas nos llevaría a un debate respecto de la ética y la gestión de los datos en contextos organizacionales.

Muchas empresas, y en particular las pertenecientes a la industria de financiera, han ido modificando, y probablemente evolucionando en su estrategia de marketing. En sus inicios las empresas concentraban su estrategia de publicidad en el uso de los medios de comunicación masivos, justamente con el objetivo de alcanzar a la mayor cantidad de usuarios posible. La efectividad de este tipo de publicidad es compleja de medir, pero según Patrutiu-Baltes (2016)ⁱ, podría concluirse que ha sido históricamente baja, debido principalmente a la gran cantidad de empresas y productos compitiendo por la atención del usuario. Muchas empresas comenzaron a utilizar nuevas estrategias de marketing basadas principalmente en datos propios de los clientes, con el objetivo de mejorar el ratio de respuesta y eficientizar las campañas de marketing. Una de las principales estrategias dentro de esta nueva rama el marketing ha sido el Telemarketing. Este consiste, básicamente, en que distintos “asesores” o “vendedores” contactan en forma directa a los potenciales clientes para ofrecer e intentar vender un determinado producto. Inicialmente ese contacto fue de forma telefónica, pero, hoy en día, gracias a la evolución de los medios de comunicación también se puede dar por otras vías o mecanismos. El telemarketing puede ser aplicado en todo tipo de industrias, pudiendo ofrecer bienes y servicios de todo tipo, pero sin duda una de las ramas que me ha aprovechado y explotado esta estrategia ha sido el sector bancario (Korkmaz, 2017)ⁱⁱ.

En principio, un banco que quisiera realizar una campaña de marketing promocionando un nuevo producto, podría contactarse de forma telefónica, con cada uno de sus actuales clientes, contándole de las supuestas ventajas de este nuevo producto, de lo conveniente que fuera para el cliente contar con el mismo, e incentivándolos a que lo adquiriera. Resulta evidente que esta no sería una estrategia para nada eficiente debido al alto costo que implicaría contactar a toda la cartera de clientes, cuando sería totalmente predecible que existe un subconjunto de cliente que a priori no tendría ningún interés en adquirir este nuevo producto en cuestión. Cabe destacar que el costo asociado a cada potencial cliente implica no solo el costo de la llamada telefónica sino el tiempo empleado por el recurso que la está realizando.

En este sentido, podemos entender que un objetivo de una empresa del sector financiero sea utilizar los datos personales propios de cada cliente para segmentar su cartera de clientes con algún propósito en particular. Considerando lo expuesto, uno de los requisitos para que una campaña de marketing sea no sólo exitosa, sino que además efectivamente eficiente es una buena segmentación del mercado. Podemos esperar que el banco tenga no solo los datos de contacto de cada uno de sus clientes, sino que muchos otros datos respecto de su situación socioeconómica, intereses y registros de gastos que les permitirían mejorar su clasificación entre cliente que se encontrarían interesados en adquirir el nuevo producto y quienes no, para concentrarse en el primer grupo. De la mano con lo afirmado por Fayyad, se requiere una nueva generación de técnicas y computacionales y herramientas para permitir la extracción de conocimiento útil de los grandes volúmenes de datos, que se encuentran en un proceso de crecimiento agigantado. Estas técnicas y herramientas son el sujeto central del campo emergente de descubrimiento de conocimiento en bases de datos (KDD) y minería de datos (Fayyad, 1996)ⁱⁱⁱ

Finalmente, a fin de lograr extraer los patrones y predecir un modelo a partir de todos los datos con los que cuenta esta organización en particular, será necesario el análisis de datos en un contexto de grandes volúmenes de Información. Esto se debe no sólo a la gran cantidad de registros con los que cuenta un banco, sino que además por la gran variedad de fuentes de donde pueden provenir dichos datos. A fin de almacenarlos y procesarlos, es necesario que un banco posea un diseño eficiente de bases de datos a fin de que los mismos puedan ser procesados de forma rápido, seguro y confiable.

El objetivo del presente trabajo consistirá entonces en la selección de un modelo de minería de datos que permite la clasificación de los clientes bancarios entre aquellos que adquirirán un determinado producto, a fin de optimizar la estrategia de Marketing Directo del Banco Santander. A tales fines, será necesario estudiar en qué consiste la metodología de extracción de conocimiento a partir de los datos (KDD) y como la misma puede ser utilizada.

Considerando lo expuesto, el objetivo general del presente trabajo consiste entonces en seleccionar el mejor modelo predictivo que permita segmentar la cartera de cliente del Banco Santander entre aquellos que estarían interesados en adquirir un determinado nuevo producto bancarios y aquellos que no. A tales fines en primer lugar será necesario también entender el concepto de Extracción de Conocimiento a partir de Datos (o Knowledge Discovery in Databases - KDD) y su aplicación a los grandes volúmenes de información en general y al caso de estudio en particular. Asimismo, identificar los distintos métodos predictivos, estudiar sus particularidades y compararlos entre sí. Calibrar cada uno de los modelos en un entorno de Python a fin de seleccionar los parámetros óptimos de cada uno de los modelos. Otro de los objetivos específicos consistirá en analizar cómo se pueden comparar la performance de cada uno de los modelos y compararlos contra un bench-mark preestablecido. Finalmente, una vez seleccionado el modelo, analizar cuáles son aquellas variables que le permitirían al banco establecer las bases de segmentación de su cartera de clientes a fin de realizar recomendaciones sobre a qué clientes direccionar su campaña de marketing.

Para analizar como un banco puede clasificar a sus clientes entre aquellos que podrían adquirir un nuevo producto y quienes no estarían interesados, mediante la aplicación de métodos analíticos descriptivos, predictivos y prescriptivos, en la sección N° 2 se comenzará por plantear el marco conceptual necesario para el desarrollo de este trabajo. Se desarrollarán los conceptos que nos permitirán comprender el rol que las técnicas de minería de datos pueden tener en la estrategia de marketing directo de una empresa. En un tercer apartado, se hará una breve descripción de la base de datos, como ha sido obtenida junto con una pequeña reflexión respecto de la responsabilidad de los datos que la misma contiene. Asimismo, en esta sección se hará una descripción de las técnicas y algoritmos planteados. En la Sección N° 4, se hará un breve análisis de los modelos propuestos y que propuestas se pueden obtener a partir de los mismos. Finalmente, en la Sección N° 5 se exponen una serie de conclusiones y trabajos finales que podrían surgir a partir de la presente investigación.

Contenido

Introducción.....	2
1. El Marketing y Minería de Datos	6
1.1. Marketing Directo	6
1.2. KDD	7
1.3. La etapa de minería de datos dentro de KDD.....	9
2. Implementación de un modelo para bancos.....	17
2.1. Base de Datos utilizada.....	17
2.2. Respecto a la privacidad de los datos	19
2.3. Estadística Descriptiva	20
2.4. Implementación de técnicas de Minería de Datos:	26
3. Selección de modelos y recomendaciones	29
Consideraciones finales y futuros trabajos	31
Anexos / Apéndices.....	34
Referencias	43

1. El Marketing y Minería de Datos

En el presente apartado se realizará un análisis de los principales conceptos teóricos que nos permitirán abordar la problemática planteada en la presente investigación.

1.1. Marketing Directo

El Marketing de Masas (“*mass marketing*”) es una de las principales ramas del marketing, y la herramienta más tradicional empleada por las empresas para promocionarse a sí mismas o sus productos. Incluye todas las herramientas y medios mediante las cuales las empresas buscan transmitir un mensaje, siendo su principal objetivo llegar a la mayor audiencia posible, mediante el uso de distintos medios de comunicación, principalmente avisos en televisión y/o radio, publicaciones en diarios y revistas, entre otros. De acuerdo con McDaniel (2009)^{iv} dentro de las ventajas identificadas del Marketing Directo se encuentran su gran masividad, y debido a esto, sus bajos costos por unidad expuesta, tanto de producción como de investigación. Sin embargo, en los últimos años, el mercado se ha puesto mucho más competitivo, y hay una variedad mucha más amplia de productos, donde todas las marcas hacen un gran esfuerzo por diferenciarse de su competencia. Si bien como se ha mencionado que este tipo de marketing llegar a audiencias altamente numerosas, dentro de sus desventajas se encuentra su baja tasa de respuesta, es decir el reducido porcentaje de personas que efectivamente compran, y más aún difícil de calcular quiénes de los clientes compraron el producto por haber visto la publicidad, volviendo así cada vez menos efectivo, y en otros términos no tan eficiente.

En contraposición al Marketing de Masas, en los últimos años se ha desarrollado una nueva rama del marketing denominada Direct Marketing (o Marketing Directo). El mismo consiste en el estudio de las características y necesidades de los potenciales clientes con el objetivo de segmentar la cartera de clientes de una empresa e identificar aquellos que potencialmente podrían y querrían adquirir un determinado producto, con el fin crear estrategias de marketing direccionadas específicamente hacia ellos con el objetivo de incentivarlos a que adquieran el producto. Algunas de estas estrategias podrían incluir el Telemarketing, mailing, couponing, televenta, email marketing, entre otras.

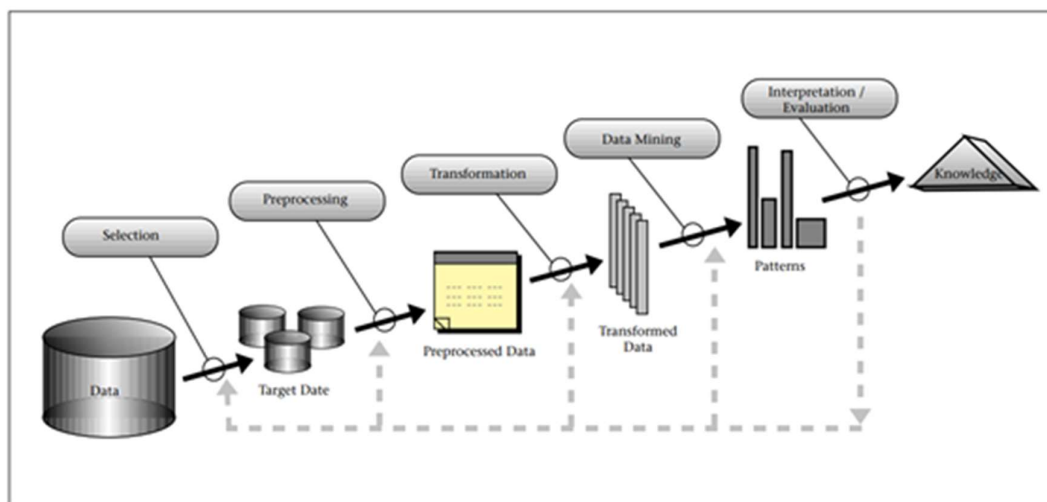
1.2. KDD

Con el transcurso del tiempo y sobre todo dependiendo del campo de estudio se le ha acuñado nombres distintos al proceso o la noción de encontrar patrones útiles a partir de los datos, incluyendo extracción de conocimiento, minería de datos, descubrimiento de información, arqueología de datos, entre tantos otros. El término Descubrimiento de conocimiento a partir de Bases de Datos (KDD) fue primeramente acuñado por Piatestsky y Schapiro^v, con el objetivo de enfatizar que el producto final del proceso radica justamente en el conocimiento con el objetivo de poder basarse en dato (data-driven). Asimismo, los mencionados autores resaltan la importancia que el proceso KDD es un concepto más amplio que las técnicas de minería de datos propiamente dichas, siendo estas tan solo un paso del proceso completo de proceso de extracción de conocimiento a partir de los datos.

De esta forma, Piatestsky y Schapiro plantean que dicho proceso podría subdividirse en 5 etapas, siendo estas no excluyentes ni exhaustivas. Según los autores, la primera etapa consiste en lo relativo a la **Preparación de los datos**, la cual tuvo un impulso justamente al estudio por las bases de datos (haciendo referencia a la segunda D del proceso KDD). Dado que muchos de los algoritmos de minería de datos suponen que los datos se encuentran cargados en la memoria, las técnicas de bases, como por ejemplo el agrupamiento y ordenamiento de datos y la optimización de consultas a la base de datos, se vuelven sumamente importante con el objetivo de ganar eficiencia en el proceso. La segunda etapa consiste en la **Selección de los datos**. Parte del proceso consiste en entender la aplicación, que se desea obtener, el conocimiento relevante ya existente y el objetivo del proceso desde la perspectiva del cliente. Una vez considerado todo esto, se debe seleccionar el set de datos con el que se va a trabajar, sean esto una recopilación de distintas fuentes de datos, o un subconjunto de toda la base existente. Comúnmente, está acompañada de la tercera etapa que describen los autores que consiste en la **Limpieza de datos**. Muchas veces esta es la tarea más dificultosa y laboriosa, y sobre la cual el cientista o especialista de datos dedica una gran proporción de su tiempo, mucho más de lo deseado. Esta etapa contempla todas las tareas necesarias para poder procesar la base de datos seleccionadas, desde la extracción de “ruido”, la estrategia para los valores perdidos, y hasta la operacionalización para que toda la información puede ser procesada por cada uno de los algoritmos (por ejemplo, ciertos algoritmos solo aceptan cierto tipo de formato en los atributos, ya sean estos valores numéricos o binarias, por lo que es necesario procesos de discretización de variables, creación de variables binarias, entre otros).

La cuarta etapa consiste propiamente en la **Minería de Datos**. Como se ha mencionado, la minería de datos es tan solo un paso en el proceso de KDD, y consiste en aplicar técnicas de análisis de datos y algoritmos de descubrimiento de patrones que, bajo limitaciones computacionales de eficiencia razonables, produzcan a una determinada enumeración de patrones o modelos sobre los datos. Sin dudas en los últimos años esta ha sido la etapa que más ha evolucionado a la par del avance en los procesadores y capacidad de las computadoras.

Finalmente, la última etapa del proceso de KDD consiste en la **Interpretación de resultados**. La última etapa consiste en la interpretación de los patrones obtenidos, y probablemente volver a algunas de las etapas anteriores para realizar una iteración, plantear nuevas variables, cambiar algún parámetro y hasta proponer un nuevo modelo. En esta etapa hay que entender que el usuario muchas veces no solo está interesado en las capacidades predictivas, sino que además encontrará mucho valor en entender qué es lo que dice el modelo.



El proceso de extracción de conocimiento a partir de datos puede ser clasificado en función de los objetivos, los cuales podemos agruparlos en dos principales. En primer lugar, el usuario o cliente puede tener el objetivo de **verificar** alguna hipótesis, sea esta ya existente o bien ya verificar que la misma se mantiene en un nuevo conjunto de datos. El segundo objetivo puede ser propiamente el de **descubrimiento**, es decir la extracción de nuevos patrones. Asimismo, el proceso de descubrimiento se puede subdividir entre *Predicción* y *Descripción*, siendo el primero utilizado para predecir el comportamiento futuro de cierta entidad, y el segundo destinado a encontrar patrones que sean entendibles para la mente humana, y una suerte de clasificación o segmentación.

1.3. La etapa de minería de datos dentro de KDD

Habiendo definido el proceso de KDD y cada una de sus etapas, se propone ahora hacer un foco en el componente de minería de datos ya que esta suele ser la que más atención recibe. Según Fayyad la minería de datos consiste en ajustar modelos a la información observada. Dicho modelo ajustado juega el rol de conocimiento inferido. El hecho de que el modelo refleje conocimiento ya sea útil, interesante o veraz por parte de los datos ya forma parte del objetivo global del proceso de KDD.

En línea con esta visión, Aggarwal (2015)^{vi} define a la minería de datos como el estudio de la recolección, limpieza, procesamiento, análisis y obtención de ideas valiosas a partir de los datos. El término de minería de datos es un amplio paraguas que es utilizado para diferentes aspectos del procesamiento de datos, y que puede ser aplicado en la resolución de una variedad de problemas sumamente amplios.

Una tercera definición de Minería de Datos puede ser la aportada por Witten, Frank y Hall (2016)^{vii}, según los cuales es el proceso de descubrir patrones a partir de los datos, donde este proceso debe realizarse de forma automática, o semiautomática. Los patrones encontrados deben ser significativos, en el sentido que deben permitir la obtención de cierta ventaja. Los patrones nos permiten realizar una predicción no trivial sobre una futura ocurrencia.

La mayoría de las técnicas de minería de datos se basan en técnicas probadas y testeadas del aprendizaje automático, reconocimiento de patrones y estadísticas. Asimismo, el verdadero modelo subyacente a la mayor parte de los algoritmos de minería de datos no es más que alguna de las siguientes: polinomios, splines, kernels y funciones básicas, entre otros modelos matemáticos

Asimismo, la etapa de Minería de Datos puede subdividirse en 3 componentes principales. El primero es la representación mediante un modelo, y consiste en el lenguaje o modelos seleccionado para descubrir los patrones. Es importante que el cientista de datos tenga un amplio entendimiento de cada una de las familias de modelos, y los supuestos que puede haber subyacente a cada uno de ellos.

El segundo componente consiste en la selección de un criterio de evaluación de los modelos. Estos criterios puede ser valor cuantitativo o una función de ajuste que mida que tan bien (o mal) un determinado modelo cumple sus objetivos en el proceso de KDD. Finalmente, el tercer componente consiste en la optimización del modelo, es decir en la selección de los parámetros que maximice el criterio de evaluación seleccionado. Una vez encontrado los parámetros óptimos, se puede volver a comenzar seleccionando una nueva familia de modelos.

Como se ha mencionado anteriormente, los dos principales objetivos de la minería de datos en la práctica suelen ser ya sea la predicción o la descripción. Aunque dicha distinción suele ser más a fines de entender los objetivos del usuario final, en la práctica los límites entre descripción y predicción no son muy estrictos, y los modelos de predicción puede ser utilizados para describir y viceversa. Ya sea con el objetivo de predicción o de descripción, la mayoría de las técnicas de minería de datos puede ser agrupadas en diversas categorías, siendo la primera de ellas los algoritmos de **Clasificación**. Weiss y Kulikowski^{viii} definieron a los modelos de clasificación como aquellos algoritmos cuyo objetivo sea aprender una función que mapee (clasifique) un registro entre una de muchas clases predefinidas. Una segunda categoría la comprenden los modelos de **Regresión**. Los modelos de regresión son aquellos que aprenden una función que mapea un registro a un cierto valor real. Asimismo, existen otros algoritmos de **Clustering**. Jain y Dubes^{ix} definieron al proceso de clusterización como una tarea descriptiva donde el objetivo es encontrar un número finito (y relativamente chico) de categorías (o clusters) que le permitan a uno lograr un mejor entendimiento del conjunto de datos. Así cada uno de los registros es asignado a un determinado cluster, o de forma más avanzada, o más de un cluster teniendo cada uno de ellos una cierta probabilidad. Finalmente, la última categoría la comprenden los algoritmos de **Sumarización**. Estos consisten en los mecanismos para lograr compactar un conjunto de registros en una única categoría. El ejemplo más simple podría ser imputar la media y el desvío de todos los registros individuales que formarían parte de un determinado subconjunto, pero existen mecanismos más sofisticados como los propuestos por Agrawal^x, a través de la derivación de reglas generales.

Habiendo realizado un desarrollo conceptual de lo que implica la extracción de conocimiento a partir de los datos, de las distintas etapas de este proceso, y hecho foco que una de las más importante es precisamente la de minería de datos, se desarrollarán a continuación algunos de los algoritmos clasificación. El espectro de posibles algoritmos es realmente amplio, y es el

investigador quien de acuerdo a su conocimiento y los datos con los que se encuentre trabajando deberá seleccionar una familia de modelos o un conjunto de familias. El o los algoritmos a utilizar depende de un conjunto de factores siendo una de las más importantes justamente el tipo de datos con los que se cuenten.

Uno de los modelos más conocidos y ampliamente utilizados son los Árboles de Decisión^{xi}. Dentro de las principales razones por las cuales este algoritmo es uno de los más populares y utilizados a la hora de realizar una clasificación se debe a que el mismo tiende a mimetizar el modo de racionalizar que usan las personas. A diferencia de otros modelos que funcionan como una caja negra, este algoritmo es fácil de entender, al menos en un caso reducido. Para crear un modelo de este tipo, se debe comenzar seleccionando un atributo que se coloca en el nodo raíz y se crea una rama para cada valor posible. Esto divide el conjunto de ejemplos en subconjuntos, uno para cada valor del atributo. El proceso continúa repitiéndose de forma recursiva para cada rama, utilizando solo aquellas instancias que realmente llegan a la rama. Si en algún momento todas las instancias en un nodo tienen la misma clasificación, deja de desarrollar esa parte del árbol. El árbol generalmente se construye como una partición jerárquica de los ejemplos de capacitación, así como el algoritmo de agrupamiento de arriba hacia abajo divide los datos jerárquicamente.

La principal diferencia de la agrupación es que el criterio de partición en el árbol de decisión se supervisa con la clase etiqueta en las instancias de entrenamiento. El algoritmo posee diversos parámetros los cuales se pueden ir modificando para optimizar la precisión de la clasificación. Dentro de estos parámetros se destacan: La cantidad de ramas (o niveles), el criterio de selección de las ramas, el tamaño mínimo para crear una nueva rama, entre otros. Asimismo, existen diversos métodos de “poda” y detención de profundidad para evitar que el modelo se sobre-ajuste.

Dentro de los modelos de clasificación otros de los modelos ampliamente reconocidos consisten en la regresión logística modela directamente las probabilidades de pertenencia a clases en términos de las variables de características con una función discriminativa. De esta forma se deben estimar los parámetros del modelo probabilístico. El análisis de regresión logística se enmarca en el conjunto de Modelos Lineales Generalizados (GLM) que usa como función de enlace la función *logit*. Las probabilidades que describen el posible resultado de un único ensayo se modelan, como una función de variables explicativas, utilizando una función logística.

$$\text{logit}(p_i) = \ln\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_k x_{k,i}.$$

De la misma manera, otro algoritmo muy utilizado a la hora de clasificación es el de Support Vector Machine (SVM)^{xii}. Este algoritmo está definido de forma natural para clasificaciones binarias, es decir cuando la variable a predecir puede tomar una de dos clases posibles, aunque también puede ser generalizado para clasificaciones de más clases. Este algoritmo se basa en un algoritmo de aprendizaje general y un núcleo de problema específico que computa el producto interno de los puntos de datos de entrada en un espacio de características. SVM realiza la clasificación al construir un hiperplano N-dimensional que separa de manera óptima los datos en dos categorías. Los modelos SVM están estrechamente relacionados con Redes neuronales. De hecho, un modelo SVM que utiliza una función de núcleo sigmoide es equivalente a una percepción de dos capas red neuronal El espacio de entrada se mapea mediante una transformación no lineal en una característica de alta dimensión espacio. El objetivo del modelado SVM es encontrar el hiperplano óptimo que separe los conjuntos de datos de tal manera que el margen entre los conjuntos de datos se maximiza. Los vectores cerca del hiperplano son los vectores de soporte. En otras palabras, el límite de decisión debe estar lo más alejado posible de los datos de ambas categorías

Las redes neuronales son un modelo de simulación del sistema nervioso humano. El sistema nervioso humano está compuesto por células, denominadas neuronas. Las neuronas biológicas están conectadas a una otro en los puntos de contacto, que se denominan sinapsis. El aprendizaje se realiza al cambiar la fuerza de las conexiones sinápticas entre las neuronas. Típicamente, la fuerza de estas conexiones cambia en respuesta a estímulos externos. Las redes neuronales pueden ser consideradas una simulación de este proceso biológico.

Como en el caso de las redes biológicas, los nodos individuales en redes neuronales artificiales se denominan neuronas. Estas neuronas son unidades de cálculo que reciben información de algunas otras neuronas, hacen cálculos en estas entradas y las alimentan en otras neuronas.

La función de cálculo en una neurona se define por los pesos en la entrada conexiones a esa neurona. Este peso puede verse como análogo a la fuerza de una conexión sináptica. Al cambiar estos pesos adecuadamente, la función de cálculo se puede aprender, lo que es análogo al aprendizaje de la fuerza sináptica en biológico Redes neuronales. El "estímulo externo" en redes neuronales artificiales para aprender estos los pesos son proporcionados por los datos de entrenamiento. La idea es modificar incrementalmente los pesos siempre que el conjunto actual de pesos haga predicciones incorrectas. Estos modelos son altamente flexibles, pero al mismo tiempo complejos, y su buen desempeño depende altamente de una correcta definición de una serie de instancias, dentro de las cuales se pueden mencionar en primer lugar la estructura de la red, es decir la cantidad de capas y nodos para cada una de ellas. En segundo lugar, en necesario definir una serie de funciones, tanto de activación como de pérdidas, y la forma de calcular los ponderadores. Asimismo, se debe definir hiper parámetros para el proceso iterativo. Las redes neuronales son un algoritmo altamente flexible, pero requieren de un alto poder computacional. La clave para la efectividad de la red neuronal es la arquitectura utilizada existiendo una gran variedad de alternativas para cada una de las instancias definidas.

A la hora de utilizar modelos para la clasificación, otra de las alternativa consiste en anidar distintos modelos, sean estos del mismo o de distinto tipo todos ellos con el objetivo de mejorar la performance del mejor modelo individual que se obtuvo hasta el momento. En caso de que se aniden modelos del mismo tipo, la idea es que cada uno de ellos tenga distintos parámetros, pero también se puede utilizar modelos pertenecientes a distintas familias de modelos. La estimación final será la ponderación de las estimaciones individuales de cada uno de los modelos. Un caso particular de los modelos ensamblados se presenta cuando los modelos que se anidan son modelos de la misma familia, comprende el caso de cuando cada modelo individual es un árbol de decisión. De esta forma se proponen un conjunto de árboles de decisión, y que juntos forman un bosque, siendo cada uno de ellos entrenados con una parte de la muestra. Cada árbol es independiente de los demás. En este caso además de definir cada uno de los parámetros presente en los árboles de decisión, se deberá definir el tamaño del bosque, o cuantos árboles se desean anidar.

El concepto fundamental detrás del bosque aleatorio es simple pero poderoso: la sabiduría de las multitudes. En ciencia de datos, la razón por la que el modelo de bosque aleatorio funciona tan bien es debido a que implementa una gran cantidad de modelos (árboles) relativamente no correlacionados que operan como comité superará a cualquiera de los modelos constituyentes individuales. La baja correlación entre modelos es la clave. Al igual que las inversiones con bajas correlaciones (como acciones y bonos) se unen para formar una cartera que es mayor que la suma de sus partes, los modelos no correlacionados pueden producir predicciones de conjunto que son más precisas que cualquiera de las predicciones individuales. La razón de este maravilloso efecto es que los árboles se protegen entre sí de sus errores individuales (siempre que no se equivoquen constantemente en la misma dirección). Si bien algunos árboles pueden estar equivocados, muchos otros árboles estarán en lo correcto, por lo que, como grupo, los árboles pueden moverse en la dirección correcta.

Al igual que el Random Forest, el clasificador Ada-boost^{xiii} es un algoritmo de clasificación en el cual se anida distintos modelos. El clasificador Ada-boost combina un algoritmo de clasificador débil para formar un clasificador fuerte. Este clasificador se basa en el hecho de que, si bien un solo algoritmo puede clasificar mal los objetos, si se combinan múltiples clasificadores con la selección del conjunto de entrenamiento en cada iteración y asignamos la cantidad correcta de peso en la votación final, se puede tener una buena puntuación de precisión para el clasificador general. Considerando esta premisa, Ada-boost, vuelve a entrenar el algoritmo de forma iterativa eligiendo el conjunto de entrenamiento basado en la precisión del entrenamiento anterior.

El peso específico de cada clasificador entrenado en cualquier iteración depende de la precisión lograda. Como se ha mencionado, Cada clasificador débil se entrena utilizando un subconjunto aleatorio del conjunto de entrenamiento general. Después de entrenar a un clasificador en cualquier nivel, Ada-boost asigna peso a cada elemento de entrenamiento. Al elemento mal clasificado se le asigna un mayor peso para que aparezca en el subconjunto de entrenamiento del siguiente clasificador con mayor probabilidad.

Asimismo, después de entrenar a cada clasificador, el peso que se le asigna a cada uno de los clasificadores también se encuentra basado en la precisión. Al clasificador más preciso se le asigna un mayor peso para que tenga más impacto en el resultado final.

Una vez seleccionados una familia de modelos o un conjunto de familias para intentar ajustar a un determinado set de datos, es necesario definir un criterio o métrica que nos permita evaluar que tan bien, o mal, se ajusta cada uno de los modelos propuestos, poder compararlos entre sí, seleccionar los parámetros óptimos para cada familia de modelos, y de todas las familias seleccionar aquel modelo que mejor nos permita explicar el set de datos propuestos. Para ello existen diversas métricas, siendo las más utilizadas las de matriz de confusión, la precisión o accuracy y el área bajo la curva ROC.

La matriz de confusión es una herramienta que permite la visualización del desempeño de un algoritmo que se emplea bajo aprendizaje supervisado. Cada columna de la matriz representa el número de predicciones de cada clase, mientras que cada fila representa a las instancias en la clase real. En caso de que se trate de una clasificación binaria, la matriz de confusión resultante será una matriz cuadrada de dos por dos.

La matriz se completa insertando en cada elemento la cantidad de observación de la clase verdadera de cada fila que son predichos bajo cada una de las clases. En el caso de las clasificaciones binarias, la matriz de confusión resultante contiene las predicciones acertadas en la diagonal principal (Verdaderos positivos y casos que no poseen el atributo y son estimados como tal). Por el otro lado, la matriz posee los casos donde el algoritmo “se equivocó”, tanto clasificando como negativamente cuando efectivamente poseía el atributo como el otro posible error de clasificar una observación positiva como negativa. Como se puede observar, las matrices de confusión es que facilitan ver si el sistema está confundiendo dos clases.

Si en los datos de entrada el número de muestras de clases diferentes cambia mucho la tasa de error del clasificador no es representativa de lo bien que realiza la tarea el clasificador. Si por ejemplo hay 990 muestras de la clase 1 y sólo 10 de la clase 2, el clasificador puede tener fácilmente un sesgo hacia la clase 1. Si el clasificador clasifica todas las muestras como clase 1 su precisión será del 99%. Esto no significa que sea un buen clasificador, pues tuvo un 100% de error en la clasificación de las muestras de la clase 2.

El *accuracy*, o la exactitud de un modelo, es una de las métricas que se pueden derivar de la matriz de confusión y nos da una idea respecto de que tan bien ajusta un modelo a un set de datos. En un caso de clasificación binaria, la exactitud se define como la proporción de casos falsos positivos y falsos negativos en relación al total de la muestra.

La curva ROC^{xiv} es una representación gráfica de la sensibilidad frente a la especificidad para un sistema clasificador binario según se varía el umbral de discriminación. En otras palabras, el área bajo la curva ROC es la representación de la razón (o ratio) de verdaderos positivos (VP = Verdaderos Positivos) en relación a los falsos positivos (FP = Falsos Positivos) también según se varía el umbral de discriminación (valor a partir del cual decidimos que un caso es un positivo). El análisis del área bajo la curva ROC, (AUC-ROC), proporciona herramientas para seleccionar los modelos posiblemente óptimos y descartar modelos. Una de las ventajas esta métrica es que la misma es independiente de la distribución de las clases en la población. El límite inferior del área bajo la curva ROC es 0.50 (lo que sería equivalente a una elección al azar), y el límite superior es 1. Cuanto mayor sea el área bajo la curva ROC mejor es el modelo no sólo en clasificar sino en ordenar a cada uno de los registros en cuanto a la confianza de su predicción.

2. Implementación de un modelo para bancos

En la presente sección se presentarán los aspectos metodológicos del presente trabajo, comenzando con lo relativo a la base de datos utilizada y como se implementaron las técnicas de minería de datos propuestas.

2.1. Base de Datos utilizada

Como se ha planteado, el objetivo del presente trabajo es desarrollar un modelo de minería de datos capaz de clasificar a los clientes de un banco entre aquellos que adquirirán un determinado producto bancario y aquellos que no. A tales fines se utilizará una base de datos que contiene registros de cada uno de los clientes del banco. Para cada uno de los clientes se poseen datos tanto personales como de otros productos que tiene contratado cada cliente, de los cuales se hace un mayor análisis en la Sección 3.2. del presente.

El presente trabajo se realizará sobre una base de datos pública del Banco Santander compartida por la misma organización a través de la plataforma Kaggle^{xv}. Kaggle Inc es una empresa perteneciente al grupo Google Inc, que se reconoce como una plataforma para el encuentro de la comunidad de científicos de datos (data scientists) e ingenieros de machine learning. La misma comenzó en el año 2017 con el objetivo de proveer una plataforma para realizar competencias de machine learning. Estas competencias son auspiciadas en la mayoría de los casos por las propias empresas dueñas de los datos quienes los ponen a disposición de la comunidad con diversos objetivos como puede ser encontrar el mejor algoritmo para el objetivo en cuestión, una fuente de reclutamiento de personal altamente capacitado para trabajar en la organización, crear cierta relación de la marca y los usuarios, entre muchos otros. Con el correr de los años, Kaggle incorporó nuevos servicios como por ejemplo los Kaggle Kernels, una sección donde los distintos usuarios pueden compartir códigos, análisis y hallazgos realizados en R o Python sobre los distintos data sets. Otro servicio ofrecido por Kaggle es una sección de educación online sobre diversos tópicos como Machine Learning, Data Visualization, Artificial Intelligence (AI), entre otros. Finalmente, otra de las secciones de Kaggle es un repositorio de base de datos públicas.

Como se comenta más adelante, el autor y responsable de la base de datos es el Propio Banco Santander, quien compartió la misma en el marco de una competencia llevada a cabo a través de la plataforma Kaggle. En la descripción de la competencia, el Organizador - Banco Santander, confirma que la base de datos se encuentra anonimizada, de forma tal de no contener información real que pueda ser asociada con ningún cliente real del banco. Asimismo, los términos y condiciones de la competencia permiten su uso siempre y cuando su uso se encuentre relacionado con leer y aprender de los datos, analizar, modificar y de forma más amplia preparar los modelos con propósito de la competencia.

De acuerdo a los términos y condiciones establecidos por Kaggle en la sección “Datasets”, apuntan no solo a ser un repositorio, sino que cada base de datos es una comunidad en donde los usuarios pueden debatir, intercambiar análisis y publicar sus ideas y códigos con el resto de los usuarios. Kaggle se encuentra preparado para admitir que las bases se encuentren en diversos formatos (como lo son CSV, JSON, SQLite, entre otros), pero incentivan a los usuarios a compartir los sets de datos en un formato abierto y accesible, con el objetivo de facilitar el trabajo y análisis de estos por toda la comunidad.

Existen dos formas en las que se pueden compartir una base de datos a través de Kaggle: La primera alternativa, y por default en Kaggle, es subir un datasets a través de un usuario de forma particular. La segunda forma es que el responsable de la base de datos sea una organización, en cuyo caso será necesario que la propia organización cree un perfil en Kaggle bajo su nombre. Esta segunda alternativa es la más conveniente en caso la base de datos a compartir haya sido creado por otro, o sea de propiedad de un tercero.

En el caso de estudio en cuestión, la base de datos utilizada formó parte de una competencia iniciada propiamente por el Banco Santander, y que una vez finalizada la base de datos quedó disponible en la sección “Datasets”.

2.2. Respeto a la privacidad de los datos

En la presente sección se reflexionará la responsabilidad del investigador en cuanto a los datos en la presente investigación bajo el marco conceptual planteado por Steinmann, Matei y Collmann^{xvi}.

Como se ha mencionado, el origen de la base de datos es la plataforma Kaggle. Kaggle Inc, siendo el autor y responsable de la base de datos el propio Banco Santander, quien compartió la misma en el marco de una competencia llevada a cabo a través de la mencionada plataforma Kaggle.

Dicha base de datos contiene cierta información sobre un conjunto de cliente del mencionado banco como, por ejemplo, País de residencia y ciudad, Sexo, Edad, Fecha de Alta, Ingreso anual, Forma de adquisición, y producto que dicho cliente tiene con el banco. Como se puede observar varios de esos datos pueden ser considerados pertenecientes al entorno privado de los propietarios de esos datos y sensibles en su utilización.

Desde la óptica del marco conceptual presentado por los mencionados autores, dicha investigación se encontraría enmarcada en un intermedio entre la reutilización de los datos y la asignación de un nuevo propósito planteado por los autores. La mencionada base de datos fue compartida por el propio Banco con el objetivo de alentar una competencia de minería de datos que permita la extracción de nuevo conocimiento por parte de la comunidad científica. En el presente caso de estudio se pretende utilizar dicha base de datos para extraer conocimiento respecto del poder de asociación y conexión de los distintos productos comercializados por el banco, que, si bien forma parte del alcance general, no es estrictamente el objetivo planteado en la competencia por el banco.

Como señalan los autores, en este caso se hace frente a dos posibles dilemas. El primero de ellos es compartido por el investigador secundarios y el propietario original de los datos (el banco), ya que este último ha compartido información sensible respecto de sus clientes. Las consecuencias del conocimiento que se pueda extraer sin dudas estarán repartidas entre estos dos actores, pero el investigador secundario no debiera excusarse a que la base de datos ha sido publicada por el banco.

Asimismo, recae sobre el investigador secundario la disyuntiva respecto de la legitimidad de la utilización de dicha base de datos con otros fines. Si bien el banco compartió la base de datos con el objetivo de que la comunidad científica extraiga conocimiento entre como se relacionan los clientes con los distintos productos que el banco ofrece, el objetivo del trabajo no es estrictamente el planteado en la competencia, pero cabe destacar que dentro de los términos y condiciones de la competencia: “se permite su uso siempre y cuando su uso se encuentre relacionado con leer y aprender de los datos, analizar, modificar y de forma más amplia preparar los modelos con propósito de la competencia”.

Por último, también existe se deben plantear los riesgos inherentes a la recombinación de los datos. Para poder compartir la base de datos de sus clientes, el banco tuvo que tomar las precauciones correspondientes para evitar que se pueda asociar e identificar a que cliente corresponde ciertos datos, ya que como se ha mencionado la base de datos publicada incluye datos sensibles. Para ello, la base de datos tuvo que pasar un proceso de anonimizarían. El riesgo aquí sería que, de alguna forma, y con el apoyo de base de datos auxiliares, se logró reconstruir la base de datos original, pudiendo por ejemplo identificar a cierta persona, sabiendo su ingreso anual, si tiene un crédito hipotecario, un auto, entre otras. De nuevo resulta importante destacar los riesgos y responsabilidades de cada parte frente a esta posibilidad, y sobre todo concientizarse que los datos no nos deberían exculparnos.

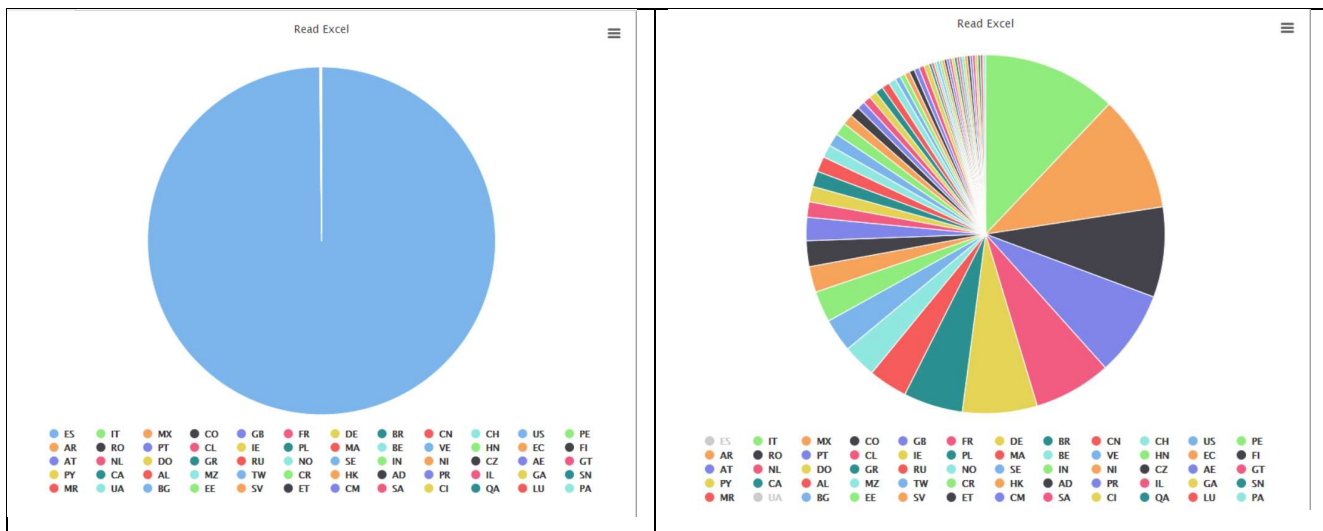
2.3. Estadística Descriptiva

Como se ha mencionado, para el desarrollo del presente trabajo se utilizó una base de datos compartida por el propio Banco Santander. La misma cuenta con más de 500.000 registros, siendo cada uno de ellos los datos relativos a cada cliente del banco. A continuación, se hace un breve resumen de las características y campos con las que cuenta la base de datos que fue utilizada para correr cada uno de los algoritmos y modelos propios de la minería de datos, sean estos propios de la base de datos original, o creados a partir de otro campo de la base de datos:

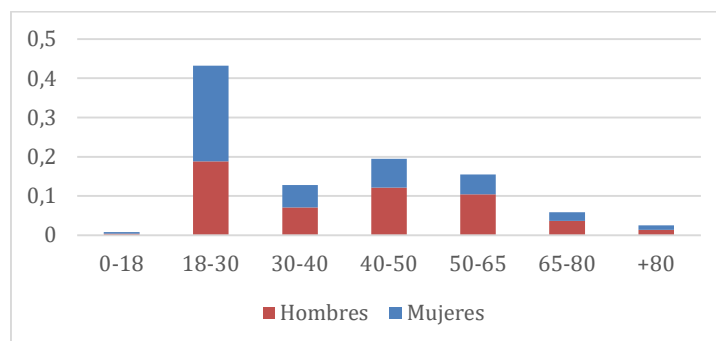
- ✓ Código del Cliente: El mismo consiste en un código para identificar al cliente dentro de esta base de datos, pero no se corresponde con el verdadero número de cliente para el banco, ya que como se ha mencionada la base de datos ha sido anonimizada por el propio banco. Cabe

destacar que no se tiene ningún otro elemento que permita identificar al cliente (como puede ser nombre, apellido, Número de documento, etc).

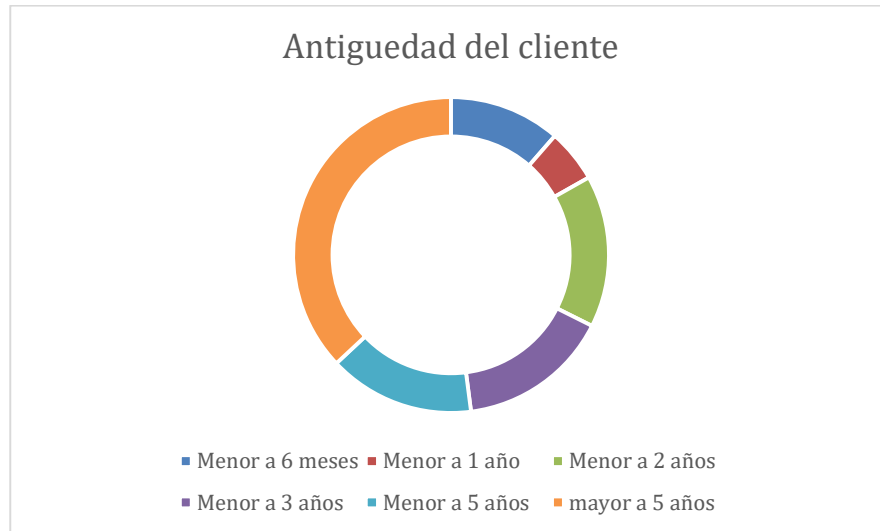
- ✓ Estado del cliente pudiendo el mismo ser Activo (A), exmpleado (B), perteneciente a una filial (F), No empleado (N) o jubilado (P).
- ✓ País de residencia. Como se puede observar la gran mayoría (más del 95%) son residentes españoles, (gráfico de la izquierda) pero si excluimos España del análisis existe una gran variedad de más de 60 alternativas distintas (gráfico de la derecha), donde las que priman son Italia, México, Colombia y Gran Bretaña.



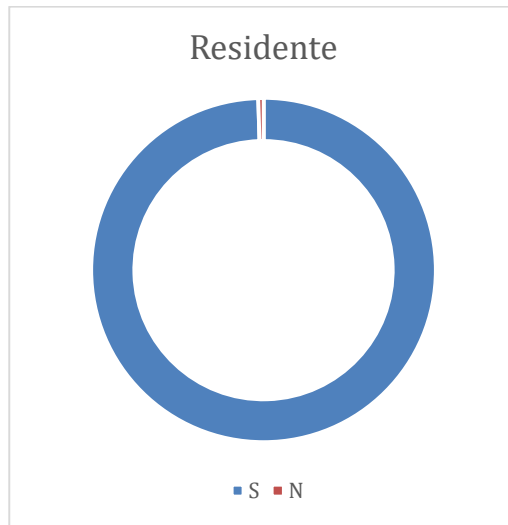
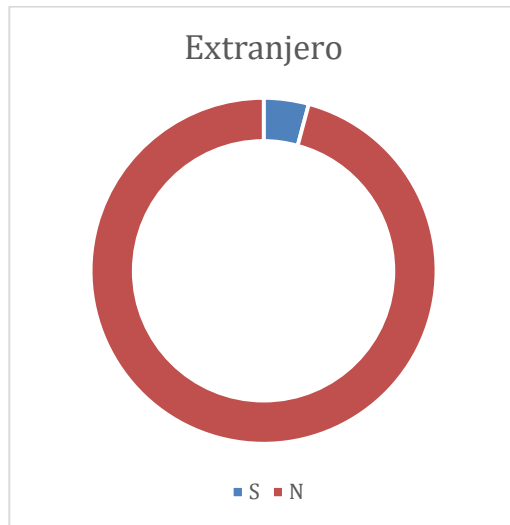
- ✓ Genero. Este atributo es binario y puede tomar los valores de 1 (hombre) o 0 (mujer).
- ✓ Edad: Refiere a la edad del cliente. El valor mínimo era 2 años y el valor máximo es 106 años, y a partir de la misma se crearon los siguientes 7 intervalos:



- ✓ Antigüedad del cliente. A partir de la fecha de alta de cada cliente, se ha calculado la antigüedad y agrupada en los siguientes intervalos

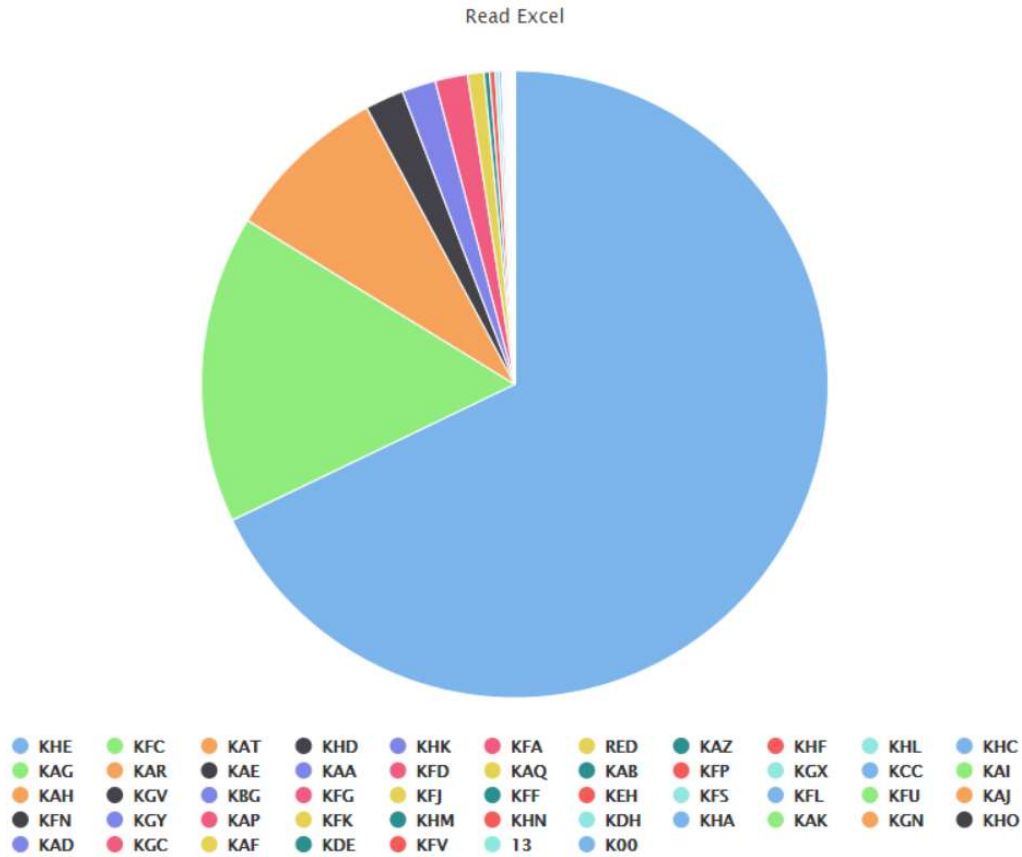


- ✓ Si el cliente es extranjero y si el cliente tiene residencia.

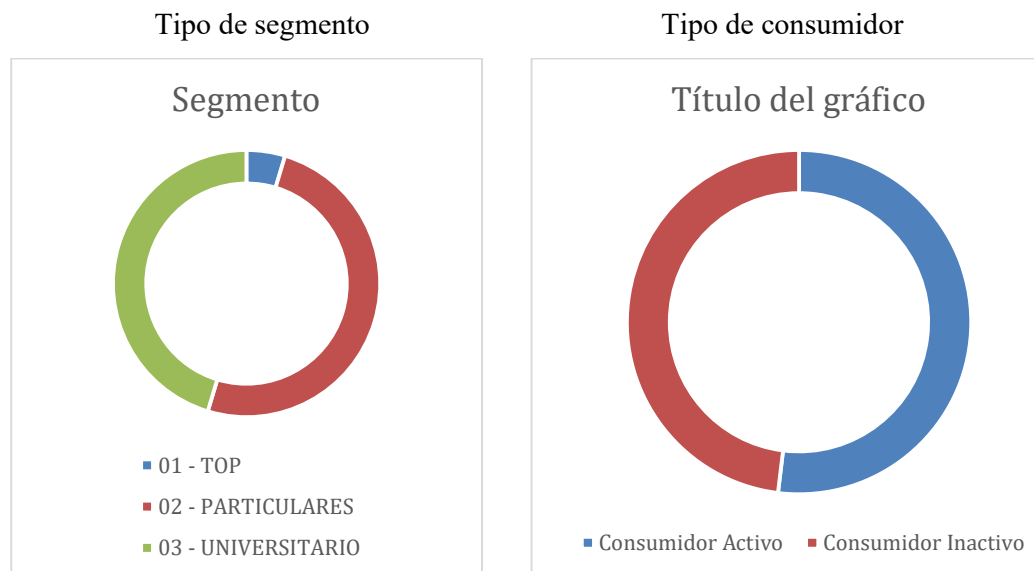


✓

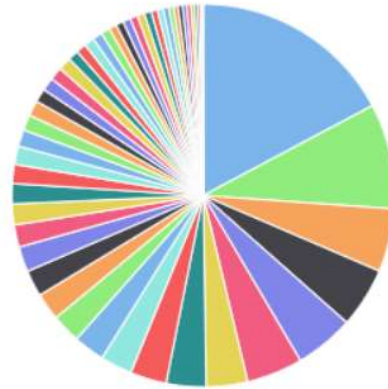
- ✓ Canal de Venta. Indica el canal de venta por el cual el cliente ha sido incorporado. Existen 43 alternativas distintas, pero como se puede observar la gran mayoría corresponde a KHE. KFC y KAT completan el podio para llegar al 90% de la muestra



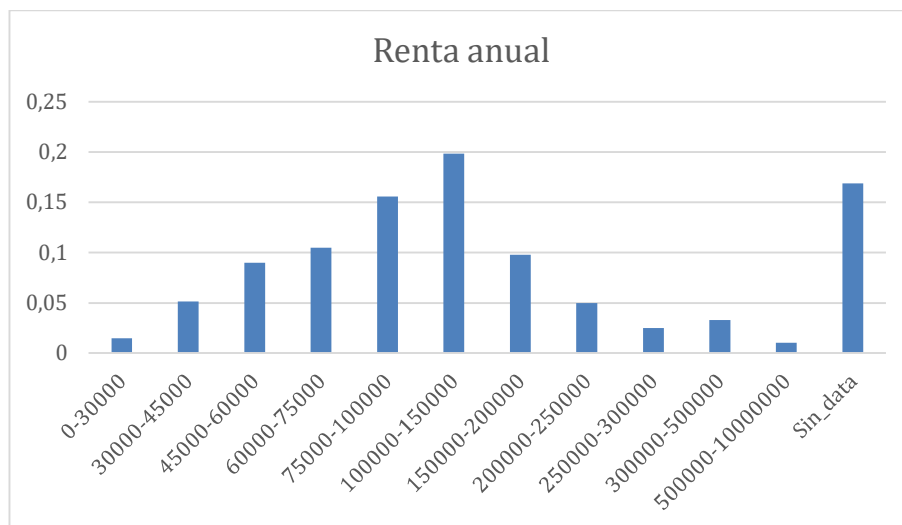
- ✓ Tipo de segmento al cual pertenece el cliente y tipo de consumidor.



✓ Domicilio declarado por el cliente



✓ Renta anual del cliente, para la cual se crearon los siguientes intervalos.



✓ Productos ya adquiridos por cada cliente.

- | | |
|-----------------------------|---------------------------|
| 1. Caja de ahorros | 12. Depósitos Largo plazo |
| 2. Garantías | 13. Hipotecas |
| 3. Cuenta Corriente | 14. Pensiones |
| 4. Cuenta Derivada | 15. Otros préstamos |
| 5. Cuenta Sueldo | 16. Impuestos |
| 6. Cuenta Juvenil | 17. Tarjetas de créditos |
| 7. Cuenta Particular Mas | 18. Securities |
| 8. Cuenta Particular | 19. Cuentas personales |
| 9. Cuenta Particular Plus | 20. Cuenta Empresas |
| 10. Depósitos corto plazo | 21. Débitos Automáticos |
| 11. Depósitos mediano plazo | 22. Otras Pensiones |
23. E-Account – Billeteras virtuales. Es la variable que se intenta predecir clasificando a los clientes ente aquellos que si los compran.



2.4. Implementación de técnicas de Minería de Datos:

Como fue desarrollado en la sección 2.3. el componente de minería de datos suele sobre el cual se pone el mayor foco en el proceso de KDD. De acuerdo a lo descrito, está puede ser subdivididas en 4 tareas, siendo la primera de ellas la definición de un conjunto de modelos que permitan clasificar a los clientes del Banco Santander entre aquellos que adquieren un determinado producto y aquellos que no, se propusieron diversos algoritmos. Como fue descrito en el apartado **2, se analizaron diversos algoritmos de clasificación**, cada modelo tiene diversos parámetros que deben ser especificados.

La segunda subetapa dentro del proceso de minería de datos consiste en un criterio de evaluación. El mismo será útil para poder comparar modelos dentro de la misma familia de modelos de forma tal de poder calibrar y seleccionar los parámetros óptimos, y asimismo el criterio de evaluación nos permitirá comparar los modelos de distinta clase. Para el presente caso de estudio se utilizará como criterio de evaluación el área bajo la curva ROC.

El tercer paso consiste justamente en calibrar cada uno de los modelos, justamente para seleccionar los parámetros óptimos. A tales fines se generó un algoritmo que permita “correr” cada uno de los modelos múltiples veces, cada una de ellas utilizando distintos parámetros, y para cada una de las corridas, se registró el área bajo la curva ROC para luego compararlas. Resulta importante destacar que con el objetivo de puntuar cada uno de los algoritmos, todos los métodos fueron evaluados utilizando el operador Cross-Validation. Este operador se utiliza a fin de evitar el sobre-ajuste (over fitting) del modelo, es decir que el modelo se sobre aprenda los datos de entrenamiento, pero luego no sea eficiente prediciendo el set de prueba. El operador consiste en subdividir toda la muestra en n submuestras, utilizando $n-1$ submuestras para el set de entrenamiento, y la submuestra restante como set de prueba, y esto lo hace de forma rotativa. Para todo este análisis se utilizó numbers of submuestras (o foldes) igual a 5. Finalmente, todos los algoritmos fueron corrido en un entorno de Python, haciéndose uso de las diferentes librerías ya desarrolladas para este fin.

El primer algoritmo en ser testeado fue el de árboles de decisión (en el Apéndice 1A se puede observar el bloque de código utilizado). Para este algoritmo se tuvo que establecer en primer lugar el criterio, que consiste en la función a ser utilizada para medir la calidad de la apertura de cada una de las ramas, siendo las opciones por elegir el *gini* y la *entropía*.

El segundo parámetro a definir es la profundidad del árbol que consiste en la cantidad de nodos con los que se quiere contar. En este caso se permitió que la misma varía entre 20 y 80, dando saltos de a 5. El tercer parámetro que se testeó fue el tamaño mínimo de muestra para que el algoritmo realice una nueva división creando un nuevo nodo. En este caso se permitió que dicho parámetro varíe entre 100 y 500, dando saltos de a 50. Como se puede observar en el apéndice 1B donde se detallan los resultados de las 208 pruebas realizadas, el mejor resultado se obtiene cuando se establece una profundidad de 20 y se divide cuando el nodo cuenta con al menos 450 observaciones, donde se obtuvo un área bajo la curva ROC de 0.856 con un error de +/- 0.053.

El segundo modelo individual en ser testeado fue el modelo de regresión logística. En este caso, el primer parámetro a definir es el algoritmo por utilizar en la optimización del problema. Asimismo, se deben definir también la función de penalización, y si se desear considerar o no un intercepto en el modelo de regresión. En el Apéndice 2A se puede observar el código utilizado, y como se puede observar en el Apéndice 2B, se obtuvieron resultados similares para todas las pruebas realizadas, siendo el mejor con AUC-ROC de 0.860 con un error de +/- 0.053.

Como se vio en el apartado de implementación de técnicas de minería de datos, una de las técnicas muy utilizadas en la minería de datos consiste en crear diferentes modelos usando muestras aleatorias con reemplazo y luego combinar o ensamblar los resultados, con el objetivo de incrementar la exactitud a partir de la combinación de las predicciones de múltiples modelos. En este sentido se utilizó el modelo de Random Forest, el cual supone la combinación de múltiples árboles de decisión individuales. De esta forma, los parámetros a definir para este modelo además de los parámetros de cada árbol, que ya han sido explicitados anteriormente, se adiciona la cantidad de árboles a combinar (`n_estimator`). En este caso para la cantidad de árboles para cada uno de los modelos se testearon desde 50 hasta 130, en saltos de a 10. En el Apéndice 3A se puede observar el bloque de código utilizado. De las más de 350 pruebas realizadas, cuyo detalle puede observar en el Apéndice 3B, la prueba con el mayor AUC-ROC se obtuvo cuando se utilizó el criterio de Gini, con 60 árboles y una máxima profundidad de 10, siendo el mismo de 0.88.

Otro de los algoritmos en ser calibrados fue el Ada-boost classifier (en el Apéndice 4A se puede observar el bloque de código utilizado). En este caso los parámetros a definir eran dos. El primero de ellos el criterio de selección, pudiendo este ser SAMME o SAMME.R. El otro

parámetro a definir es la cantidad de estimadores, para lo cual se testearon desde 20 hasta 100 en intervalos de a 100.

Como se puede observar en el apéndice 4B donde se detallan los resultados de las 16 pruebas realizadas, el mejor resultado se obtiene cuando se utiliza el algoritmo SAMME.R y una cantidad de 50 estimadores, lográndose un área bajo la curva ROC, levemente superior del 0.867, pero además lográndose disminuir el error a +/- 0.049.

Uno de los algoritmos de minería de datos con crecimiento más vertiginoso en los últimos tiempos ha sido las redes neuronales, principalmente por su gran flexibilidad y versatilidad para adaptarse ante distintos conjuntos de datos. Esta versatilidad se debe principalmente a la gran cantidad de parámetros y estructuras con las que se puede definir y planear una red neuronal. En este caso, el primer parámetro que se tuvo que definir es el algoritmo que el modelo va a utilizar para entrenarse, pudiéndose seleccionar entre Stochastic Gradient Descent, Adam o L-BFGS. SGD (Descenso de gradiente estocástico) actualiza los parámetros utilizando el gradiente de la función de pérdida con respecto a un parámetro que necesita adaptación, es decir donde es la tasa de aprendizaje que controla el tamaño del paso en la búsqueda de espacio de parámetros. es la función de pérdida utilizada para la red. Adam es similar al SGD en el sentido de que es un optimizador estocástico, pero puede ajustar automáticamente la cantidad para actualizar los parámetros en función de las estimaciones adaptativas de los momentos de orden inferior. L-BFGS es un solucionador que se aproxima a la matriz de Hessiana que representa la derivada parcial de segundo orden de una función. Además, se aproxima el inverso de la matriz de Hesse para realizar actualizaciones de parámetros. La implementación utiliza la versión Scipy de L-BFGS.

El siguiente parámetro a definir es la función de activación, siendo esta la función que permite que la información fluya de una serie de nodos a la siguiente. En este caso se testearon las funciones de identidad, la logística, la “tanh” y “relu”. Finalmente es necesario definir la estructura de perceptrones de cada una de las capas de la red. En este caso, se propusieron diversos esquemas con distintas cantidades de capas. En el Apéndice 5A se puede observar un detalle de todas las pruebas realizadas y los resultados obtenidos.

3. Selección de modelos y recomendaciones

Como se ha detallado en el apartado anterior, se han analizado, testeado y optimizados diversos métodos analíticos predictivos, obteniéndose para cada uno de ellos una determinada área bajo la curva ROC. A continuación, se presenta un cuadro donde se resumen todos los resultados obtenidos.

Algoritmo	#	Criterio/ Algoritmo	n_estimator (cantidad de árboles)	Maxima Profundiad	min_sample s_split	AUC	Desvío_AUC
Decision Tree	8	gini		20	450	0.85634913	0.05376201
Regresión Logística	1	newton-cg	12	FALSE		0.86065181	0.05352686
Random Forest	123	gini	60	10	200	0.88025245	0.039946
AdaBoost	4	SAMME	50			0.86704784	0.04913116
Redes Neuronales	16	lbfgs	tanh	(11,9,7,5,3)		0.88782224	0.04312936

Si bien todos los modelos han performado bastante similar con área bajo la curva ROC (AUC-ROC) relativamente altas para todas las simulaciones, se puede observar que, el modelo de Redes Neuronales con la estructura y parámetros definidos es el que mejor AUC-ROC arroja. No obstante, el AUC-ROC obtenidos con el modelo de Random Forest, si bien es levemente inferior, tiene la ventaja de haber arrojado un menor error, es decir menor variación entre cada uno de los resultados cuando se realizó el Cross Valiadation, lo que da la idea de que se trata de un modelo más estable, y que probablemente pueda tener una mejor performance ante un set de datos distintos al que fue utilizado para entrenar.

Otra de las razones que podría inclinar la balanza en favor del Modelo de Random Forest, es que el mismo se trata de un modelo más simple y fácil de explicar a un tercero, ya que el mismo se basa em reglas claras (cada una de las ramas de decisión del árbol). En contraposición, una red neuronal suele ser un modelo no solo más complejo, sino que también funciona como una caja negra. De esta forma, es posible extraer del modelo de Random Forest cuales son las variables relevantes en la clasificación (las ramas principales), lo cual es de suma importancia para la organización, ya que podrían funcionar como factores de segmentación con el objetivo de direccionar sus campañas y esfuerzos de marketing.

Con el bloque de código que se expone en el Apéndice 6, se puede extraer las variables mas relevantes del modelo seleccionado. De esta forma se puede concluir que aquellos clientes que tienden a comprar el producto objeto de estudio del presente trabajo tienden a tener las siguientes características:

- ✓ El cliente ya posee una Cuenta corrientes con la organización
- ✓ El cliente posee una antigüedad mayor a los 3 años.
- ✓ El cliente ya posee una cuenta junior (o junior plus) con la organización
- ✓ Tienen como canal de venta KHE o KHD
- ✓ Se trata de un cliente particular (no empresa)
- ✓ Clientes menores de 40 años (frangas etarias 18-30 y 30-40).
- ✓ El cliente tiene una renta menor a 75.000.

Consideraciones finales y futuros trabajos

Durante el presente trabajo se estudió como en los últimos años, todas las organizaciones, pero en particular las de la industria financiera, han comenzado a recopilar datos de sus clientes de forma sistemática. Asimismo, se analizó la importancia que pueden tener esos datos a la hora de querer clasificar la cartera de clientes de una organización, y los beneficios de una buena segmentación de mercados y potenciales usuarios con el fin de mejorar la estrategia de marketing directo de una organización.

En los últimos años se ha evidenciado un crecimiento constante y vertiginoso de las herramientas informáticas que han permitido avanzar no solo en el almacenamiento de grandes volúmenes de información, sino que además en su procesamiento a fin de poder analizar los mismos y generar modelos que ayuden a la toma de decisiones. Como se expuso en la sección 2.1. resulta de vital importancia dirigir y segmentar el público al cual se destina toda estrategia de marketing, no solo a fines de reducir los costos que la mismas conllevan, sino que también reducir el destino de recursos (por ejemplo, de telemarketers realizando las llamadas) sino que también del costo reputacional que tiene “la molestia” en contactar a alguien que no desea adquirir el servicio.

Asimismo, durante el desarrollo del presente trabajo se describió el proceso de extracción de conocimiento a partir de los datos, realizando una descripción de cada una de las etapas. Se examinó la importancia de la recolección de los datos, y aún mas importante el preprocesamiento, la limpieza y la preparación de estos para luego poder ser procesados por los distintos métodos analíticos predictivos. En particular, se hizo foco en la etapa del proceso de KDD que más desarrollo ha tenido en los últimos años referente a las técnicas de minería de datos propiamente dichos, analizándose los algoritmos de clasificación más utilizados y realizándose una breve reseña de sus particularidades mas importantes y ventajas de cada uno.

El objetivo particular del presente trabajo fue aplicar los mencionados modelos al caso de estudio en particular, con el objetivo de lograr el mejor modelo que pudiera clasificar a los clientes del banco y predecir quienes adquirirían un determinado producto. Antes de poder calibrar cada uno de los modelos, fue necesario detallar el origen de la base de datos, realizar una descripción del set de datos original y del proceso de preparación, con un detalle de las variables que fueron creadas a partir de los datos existentes.

Debido al origen de la base de datos, no menos importante en la sección 3.1.3 se realizó una reflexión de la investigación bajo el marco conceptual planteado por Steinmann, Matei y Collmann y de la corresponsabilidad del investigador debido a utilizar un set de datos público de internet.

En el apartado **2. Implementación de técnicas de Minería de Datos** se detalló como se implementaron los distintos modelos en un entorno de Python, haciéndose un detalle de los diversos algoritmos utilizados, que parámetros se testearon para cada uno de los modelos y los resultados obtenidos. En el apartado N° 3 no solo se hizo un resumen de todos los resultados, sino que se concluyó que el mejor modelo para clasificar a los clientes de esta organización en relación al nuevo producto es el modelo de Random Forest. Mas importante aún en dicho apartado, se pudo extraer cuales son las características más relevantes a la hora de clasificar a los clientes, para que la institución financiera las tenga en cuenta a la hora de segmentar su cartera de clientes y/o potenciales usuario con el fin direccionar hacia ellos los esfuerzos y recursos de las campañas de marketing directo a realizar.

Por concluir el presente trabajo de investigación, quedan pendientes algunas consideraciones para ser abordadas en próximos trabajos en un futuro cercano, que bien podrían ser desarrolladas por este investigador, o bien por cualquier otro entusiasta de la extracción del conocimiento a partir de los datos que podrían tomar este trabajo como punto de partida.

En primer lugar, se podrían evaluar implementar alguno de los modelos desarrollados en el presente trabajo en un entorno de producción. Esto implicaría tomar un modelo ya optimizado en sus parámetros y moverlo hacia una plataforma que permita la interacción directa con el usuario. De esta forma, le permitiría a la organización segmentar una determinada base de datos de clientes, ordenando los registros de acuerdo a la probabilidad que el modelo le estima a cada usuario de éxito en la adquisición del producto de interés. Las ventajas en la implementación de tal modelo para la organización son diversas, destacándose principalmente el ordenamiento de los usuarios, el descarte de cierto grupo de clientes que se espera no esté interesados en adquirir el producto, con la posibilidad de redireccionar dichos recursos y esfuerzos hacia los clientes que si se espera se encuentren interesados en el producto.

Por último, otra de los trabajos que podrían realizarse como complemento de todo lo desarrollado en el presente es un análisis de asociación entre los diversos productos. Este análisis permitiría estudiar y extraer conocimiento respecto de como los diversos usuarios del banco se comportan en relación a los mas de 25 productos que la institución financiera ofrece, con el objetivo de determinar, en caso de existir, que productos suelen adquirirse de forma simultánea, y cuales podrían llegar a ser sustitutos. Asimismo, todo este análisis permitiría potenciar los modelos desarrollados de forma de lograr un sistema de recomendación completo. Este nuevo modelo permitiría, no solo clasificar si un determinado cliente estaría interesado o no en el nuevo producto, sino que además permitiría predecir, dado el conjunto de productos que cada cliente ya tenga, cual o cuales serían los productos adicionales que el cliente podría estar interesado en adquirir.

Anexos / Apéndices

Apéndice 1A – Código de Decision Tree

```
## 0. Importar los paquetes necesarios junto con el set de datos

import pandas as pd
import sklearn as sk

from google.colab import files
uploaded = files.upload()

## 1. LEER DATOS...

data = pd.read_excel("final_v2.xlsx")

# 2. PREPARAR DATOS
#data = pd.get_dummies(data, columns=['area', 'salario'])

X = data[data.columns.drop(['ncodpers', 'ind_ecue_fin_ult1']).values]
y = data['ind_ecue_fin_ult1'].values

## DECISION TREE
# 2. VALIDACION CRUZADA EN 5 PARTES

from sklearn import model_selection
from sklearn.tree import DecisionTreeClassifier
from sklearn import tree
from sklearn.model_selection import cross_val_score

for criterion in ("gini", "entropy"):
    for max_depth in range(20, 85, 5):
        for min_samples_split in range(100, 500, 50):
            clf = tree.DecisionTreeClassifier(criterion=criterion, max_depth=max_depth, min_samples_split=min_samples_split)
            scores = model_selection.cross_val_score(clf, X, y, scoring="roc_auc", cv=5)
            print(criterion, max_depth, min_samples_split, "exactitud:", scores.mean(), scores.std(), '|')
```

Apéndice 1B – Resultados obtenidos de Decision Tree

Algoritmo	#	Criterio	Maxima Profundiad	min_samples_split	AUC	Desvío_AUC	Algoritmo	#	Criterio	Maxima Profundiad	min_samples_split	AUC	Desvío_AUC
Decision Tree	1	gini	20	100	0.84541997	0.05539557	Decision Tree	105	entropy	20	100	0.84492109	0.0607289
Decision Tree	2	gini	20	150	0.84827934	0.05627461	Decision Tree	106	entropy	20	150	0.84931736	0.05938861
Decision Tree	3	gini	20	200	0.85165315	0.05408099	Decision Tree	107	entropy	20	200	0.85111831	0.06037991
Decision Tree	4	gini	20	250	0.85209538	0.05498287	Decision Tree	108	entropy	20	250	0.85205111	0.06019875
Decision Tree	5	gini	20	300	0.85319558	0.0554698	Decision Tree	109	entropy	20	300	0.85418945	0.05977929
Decision Tree	6	gini	20	350	0.85510611	0.05450246	Decision Tree	110	entropy	20	350	0.85503972	0.06027763
Decision Tree	7	gini	20	400	0.85590686	0.05377214	Decision Tree	111	entropy	20	400	0.8555802	0.06059681
Decision Tree	8	gini	20	450	0.85634913	0.05376201	Decision Tree	112	entropy	20	450	0.85595578	0.06115912
Decision Tree	9	gini	25	100	0.83816077	0.06052033	Decision Tree	113	entropy	25	100	0.83734121	0.06364194
Decision Tree	10	gini	25	150	0.84359916	0.06052422	Decision Tree	114	entropy	25	150	0.8435574	0.06388655
Decision Tree	11	gini	25	200	0.84765377	0.05834898	Decision Tree	115	entropy	25	200	0.8459329	0.06460686
Decision Tree	12	gini	25	250	0.84905353	0.05947605	Decision Tree	116	entropy	25	250	0.84784486	0.06467099
Decision Tree	13	gini	25	300	0.85120456	0.05932092	Decision Tree	117	entropy	25	300	0.85027236	0.06354401
Decision Tree	14	gini	25	350	0.85286598	0.05878263	Decision Tree	118	entropy	25	350	0.8515961	0.06419189
Decision Tree	15	gini	25	400	0.85364911	0.05819853	Decision Tree	119	entropy	25	400	0.85259606	0.06448391
Decision Tree	16	gini	25	450	0.8545081	0.05833948	Decision Tree	120	entropy	25	450	0.85284219	0.06532672
Decision Tree	17	gini	30	100	0.83395148	0.05948614	Decision Tree	121	entropy	30	100	0.83349903	0.06191058
Decision Tree	18	gini	30	150	0.84112973	0.06009055	Decision Tree	122	entropy	30	150	0.84084321	0.0624113
Decision Tree	19	gini	30	200	0.84570753	0.05808925	Decision Tree	123	entropy	30	200	0.84368236	0.0635035
Decision Tree	20	gini	30	250	0.84741931	0.0590517	Decision Tree	124	entropy	30	250	0.8464424	0.06347225
Decision Tree	21	gini	30	300	0.84975645	0.05914737	Decision Tree	125	entropy	30	300	0.84892564	0.06264123
Decision Tree	22	gini	30	350	0.85174858	0.05852214	Decision Tree	126	entropy	30	350	0.8504484	0.06317042
Decision Tree	23	gini	30	400	0.85277775	0.05795745	Decision Tree	127	entropy	30	400	0.85144394	0.0638651
Decision Tree	24	gini	30	450	0.85373856	0.05806859	Decision Tree	128	entropy	30	450	0.85171384	0.0647986
Decision Tree	25	gini	35	100	0.83171566	0.05948957	Decision Tree	129	entropy	35	100	0.83126926	0.06174583
Decision Tree	26	gini	35	150	0.83973395	0.05997108	Decision Tree	130	entropy	35	150	0.83917362	0.06241283
Decision Tree	27	gini	35	200	0.84476152	0.05796702	Decision Tree	131	entropy	35	200	0.84296757	0.06307811
Decision Tree	28	gini	35	250	0.84674322	0.05907097	Decision Tree	132	entropy	35	250	0.84535558	0.06291449
Decision Tree	29	gini	35	300	0.84892988	0.05952631	Decision Tree	133	entropy	35	300	0.84847314	0.06208145
Decision Tree	30	gini	35	350	0.85133969	0.05867585	Decision Tree	134	entropy	35	350	0.84983339	0.06280612
Decision Tree	31	gini	35	400	0.85264708	0.05785859	Decision Tree	135	entropy	35	400	0.85113601	0.06331388
Decision Tree	32	gini	35	450	0.85345183	0.05797771	Decision Tree	136	entropy	35	450	0.85137998	0.06452761
Decision Tree	33	gini	40	100	0.83092743	0.05917786	Decision Tree	137	entropy	40	100	0.8301915	0.06171159
Decision Tree	34	gini	40	150	0.83904971	0.0595955	Decision Tree	138	entropy	40	150	0.83836963	0.06225636
Decision Tree	35	gini	40	200	0.84439935	0.05779105	Decision Tree	139	entropy	40	200	0.84229409	0.06349689
Decision Tree	36	gini	40	250	0.84658675	0.05885436	Decision Tree	140	entropy	40	250	0.84527485	0.06314514
Decision Tree	37	gini	40	300	0.84909203	0.05890632	Decision Tree	141	entropy	40	300	0.8483453	0.06202388
Decision Tree	38	gini	40	350	0.85116185	0.05847665	Decision Tree	142	entropy	40	350	0.84982522	0.06269998
Decision Tree	39	gini	40	400	0.85249103	0.05784328	Decision Tree	143	entropy	40	400	0.85099992	0.0634754
Decision Tree	40	gini	40	450	0.853399	0.05786136	Decision Tree	144	entropy	40	450	0.85120951	0.06454456
Decision Tree	41	gini	45	100	0.83066242	0.05907057	Decision Tree	145	entropy	45	100	0.82960015	0.06168362
Decision Tree	42	gini	45	150	0.83907208	0.05953261	Decision Tree	146	entropy	45	150	0.83850453	0.06208412
Decision Tree	43	gini	45	200	0.8444495	0.05777526	Decision Tree	147	entropy	45	200	0.84196289	0.06338926
Decision Tree	44	gini	45	250	0.84649289	0.05892083	Decision Tree	148	entropy	45	250	0.84539456	0.0631916
Decision Tree	45	gini	45	300	0.84904672	0.05896623	Decision Tree	149	entropy	45	300	0.84855931	0.06210911
Decision Tree	46	gini	45	350	0.85114054	0.05846485	Decision Tree	150	entropy	45	350	0.84974137	0.06289892
Decision Tree	47	gini	45	400	0.85247743	0.05787267	Decision Tree	151	entropy	45	400	0.85100747	0.0633468
Decision Tree	48	gini	45	450	0.85344354	0.05796655	Decision Tree	152	entropy	45	450	0.85135037	0.0643848
Decision Tree	49	gini	50	100	0.83085276	0.05901489	Decision Tree	153	entropy	50	100	0.82979837	0.06166255
Decision Tree	50	gini	50	150	0.8390389	0.05960778	Decision Tree	154	entropy	50	150	0.83818361	0.06254035
Decision Tree	51	gini	50	200	0.8444007	0.05769651	Decision Tree	155	entropy	50	200	0.84249615	0.06303274
Decision Tree	52	gini	50	250	0.84652937	0.05890227	Decision Tree	156	entropy	50	250	0.84489277	0.06318304
Decision Tree	53	gini	50	300	0.84906846	0.05886234	Decision Tree	157	entropy	50	300	0.8483232	0.06250003
Decision Tree	54	gini	50	350	0.85115354	0.05845613	Decision Tree	158	entropy	50	350	0.84979176	0.06281108
Decision Tree	55	gini	50	400	0.85256438	0.05780668	Decision Tree	159	entropy	50	400	0.85089175	0.06359146
Decision Tree	56	gini	50	450	0.85344879	0.0578688	Decision Tree	160	entropy	50	450	0.85127453	0.06457207
Decision Tree	57	gini	55	100	0.83073036	0.05891246	Decision Tree	161	entropy	55	100	0.82990114	0.06151197
Decision Tree	58	gini	55	150	0.83904983	0.05955572	Decision Tree	162	entropy	55	150	0.83824467	0.06212303
Decision Tree	59	gini	55	200	0.84439218	0.05782532	Decision Tree	163	entropy	55	200	0.84244929	0.06292459
Decision Tree	60	gini	55	250	0.84656487	0.05882674	Decision Tree	164	entropy	55	250	0.84503414	0.06310641
Decision Tree	61	gini	55	300	0.84904639	0.0588553	Decision Tree	165	entropy	55	300	0.84868935	0.06194403
Decision Tree	62	gini	55	350	0.85115725	0.05854315	Decision Tree	166	entropy	55	350	0.84993204	0.06252659

Apéndice 2A – Código de Regresión Logística.

```
## 0. Importar los paquetes necesarios junto con el set de datos

import pandas as pd
import sklearn as sk

from google.colab import files
uploaded = files.upload()

## 1. LEER DATOS...

data = pd.read_excel("final_v2.xlsx")

# 2. PREPARAR DATOS
#data = pd.get_dummies(data, columns=['area', 'salario'])

X = data[data.columns.drop(['ncodpers', 'ind_ecue_fin_ult1'])].values
y = data['ind_ecue_fin_ult1'].values

from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score

for solversrt in ("newton-cg", "lbfgs", "sag", "saga"):
    for penalty in ("l1", "l2", "elasticnet"):
        for fit_intercept in ("TRUE", "FALSE"):
            clf = LogisticRegression(penalty=penalty, fit_intercept=fit_intercept)
            scores = model_selection.cross_val_score(clf, X, y, scoring="roc_auc", cv=5)
            print(solversrt, penalty, fit_intercept, "exactitud:", scores.mean(), scores.std(), '|')
```

Apéndice 2B – Resultados obtenidos de Regresión Logística

Algoritmo	#	solversrt	penalty	intercept	AUC	Desvío_AUC
Regresión Logística	1	newton-cg	l2	FALSE	0.86065181	0.05352686
Regresión Logística	2	newton-cg	l2	TRUE	0.86065181	0.05352686
Regresión Logística	3	newton-cg	l1	FALSE	0.86011496	0.05399378
Regresión Logística	4	newton-cg	l1	TRUE	0.86012104	0.05399649

Apéndice 3A – Código de Random Forest

```
## 0. Importar los paquetes necesarios junto con el set de datos

import pandas as pd
import sklearn as sk

from google.colab import files
uploaded = files.upload()

## 1. LEER DATOS...

data = pd.read_excel("final_v2.xlsx")

# 2. PREPARAR DATOS
#data = pd.get_dummies(data, columns=['area', 'salario'])

X = data[data.columns.drop(['ncodpers', 'ind_ecue_fin_ult1'])].values
y = data['ind_ecue_fin_ult1'].values

## DECISION TREE
# 2. VALIDACION CRUZADA EN 5 PARTES

from sklearn import model_selection
from sklearn import ensemble
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import RandomForestClassifier

for criterion in ("gini", "entropy"):
    for n_estimators in range(50, 130, 10):
        for max_depth in range(20, 85, 5):
            for min_samples_split in range(100, 500, 50):
                clf = RandomForestClassifier(criterion=criterion, n_estimators=
n_estimators, max_depth=max_depth, min_samples_split=min_samples_split)
                scores = model_selection.cross_val_score(clf, X, y, scoring="roc_auc", cv=5)
                print(criterion, n_estimators, max_depth, min_samples_split, "exactitud:", scores.mean(), scores.std(), '|')
```


Apéndice 4A – Código de AdaBoost

```
## 0. Importar los paquetes necesarios junto con el set de datos

import pandas as pd
import sklearn as sk

from google.colab import files
uploaded = files.upload()

## 1. LEER DATOS...

data = pd.read_excel("final_v2.xlsx")

# 2. PREPARAR DATOS
#data = pd.get_dummies(data, columns=['area', 'salario'])

X = data[data.columns.drop(['ncodpers', 'ind_ecue_fin_ult1'])].values
y = data['ind_ecue_fin_ult1'].values

from sklearn import model_selection
from sklearn.ensemble import AdaBoostClassifier

for algorithm in ("SAMME", "SAMME.R"):
    for n_estimators in range(20,100,10):
        clf = sk.ensemble.AdaBoostClassifier(n_estimators=n_estimators)
        scores = sk.model_selection.cross_val_score(clf, X, y, scoring="roc_
_auc", cv=5)
        print(algorithm, n_estimators, "exactitud:", scores.mean(), scores.
std(), '|')
```

Apéndice 4A – Resultados obtenidos de AdaBoost

Algoritmo	#	Algorithm	n parameters	AUC	Desvío_AUC
AdaBoost	1	SAMME	20	0.86266812	0.05091928
AdaBoost	2	SAMME	30	0.86333534	0.04994779
AdaBoost	3	SAMME	40	0.8641961	0.05058345
AdaBoost	4	SAMME	50	0.86704784	0.04913116
AdaBoost	5	SAMME	60	0.86681234	0.04948466
AdaBoost	6	SAMME	70	0.86579564	0.04934866
AdaBoost	7	SAMME	80	0.86578128	0.04915402
AdaBoost	8	SAMME	90	0.86584354	0.04915836
AdaBoost	9	SAMME.R	20	0.86266812	0.05091928
AdaBoost	10	SAMME.R	30	0.86333534	0.04994779
AdaBoost	11	SAMME.R	40	0.8641961	0.05058345
AdaBoost	12	SAMME.R	50	0.86704784	0.04913116
AdaBoost	13	SAMME.R	60	0.86681234	0.04948466
AdaBoost	14	SAMME.R	70	0.86579564	0.04934866
AdaBoost	15	SAMME.R	80	0.86578128	0.04915402
AdaBoost	16	SAMME.R	90	0.86584354	0.04915836

Apéndice 5A – Código de Redes Neuronales

```

from sklearn import model_selection
from sklearn import ensemble
from sklearn.model_selection import cross_val_score
from sklearn.neural_network import MLPClassifier

for solver in ("lbfgs", "sgd", "adam"):
    for activation in ("identity", "logistic", "tanh", "relu"):
        for hidden_layer_sizes in ((5,3), (7,5,3), (9,7,5,3), (11,9,7,5,3),
            (13,11,9,7,5,3), (21,15,9,3)):
            clf = MLPClassifier(solver=solver, activation=activation, alpha=1
                e-5, hidden_layer_sizes=hidden_layer_sizes, random_state=1)
            scores = sk.model_selection.cross_val_score(clf, X, y, scoring="r
                oc_auc", cv=5)
            print(solver, activation, hidden_layer_sizes, scores.mean(), scor
                es.std(), '|')
    
```

Apéndice 5B – Resultados obtenidos de Redes Neuronales

Algoritmo	#	Algoritmo	Funcion de activacion	Estructura	AUC	Desvío_AUC	Algoritmo	#	Algoritmo	Funcion de activacion	Estructura	AUC	Desvío_AUC
Redes Neuronales	1	lbfgs	identity	(5,3)	0.85811557	0.0578678	Redes Neuronales	37	sgd	tanh	(5,3)	0.88156738	0.04614304
Redes Neuronales	2	lbfgs	identity	(7,5,3)	0.85928944	0.05501596	Redes Neuronales	38	sgd	tanh	(7,5,3)	0.87466926	0.05049405
Redes Neuronales	3	lbfgs	identity	(9,7,5,3)	0.85701674	0.05775205	Redes Neuronales	39	sgd	tanh	(9,7,5,3)	0.8831407	0.04431997
Redes Neuronales	4	lbfgs	identity	(11,9,7,5,3)	0.85960432	0.05523463	Redes Neuronales	40	sgd	tanh	(11,9,7,5,3)	0.88636741	0.03915643
Redes Neuronales	5	lbfgs	identity	(13,11,9,7,5,3)	0.8595505	0.05570637	Redes Neuronales	41	sgd	tanh	(13,11,9,7,5,3)	0.88414322	0.04041246
Redes Neuronales	6	lbfgs	identity	(21,15,9,3)	0.85963468	0.05621495	Redes Neuronales	42	sgd	tanh	(21,15,9,3)	0.88427476	0.0442243
Redes Neuronales	7	lbfgs	logistic	(5,3)	0.87682343	0.0470888	Redes Neuronales	43	sgd	relu	(5,3)	0.86674537	0.04934428
Redes Neuronales	8	lbfgs	logistic	(7,5,3)	0.87278381	0.04775394	Redes Neuronales	44	sgd	relu	(7,5,3)	0.87939539	0.04651211
Redes Neuronales	9	lbfgs	logistic	(9,7,5,3)	0.69669365	0.08153834	Redes Neuronales	45	sgd	relu	(9,7,5,3)	0.87996558	0.04120093
Redes Neuronales	10	lbfgs	logistic	(11,9,7,5,3)	0.44824952	0.06788125	Redes Neuronales	46	sgd	relu	(11,9,7,5,3)	0.87501784	0.04277814
Redes Neuronales	11	lbfgs	logistic	(13,11,9,7,5,3)	0.3564934	0.14688834	Redes Neuronales	47	sgd	relu	(13,11,9,7,5,3)	0.877061	0.04895362
Redes Neuronales	12	lbfgs	logistic	(21,15,9,3)	0.27859535	0.10235737	Redes Neuronales	48	sgd	relu	(21,15,9,3)	0.88413117	0.04128863
Redes Neuronales	13	lbfgs	tanh	(5,3)	0.87609877	0.04437908	Redes Neuronales	49	adam	identity	(5,3)	0.85760957	0.05769905
Redes Neuronales	14	lbfgs	tanh	(7,5,3)	0.88414944	0.04536765	Redes Neuronales	50	adam	identity	(7,5,3)	0.86202298	0.05307184
Redes Neuronales	15	lbfgs	tanh	(9,7,5,3)	0.88055858	0.04949618	Redes Neuronales	51	adam	identity	(9,7,5,3)	0.86025403	0.05271358
Redes Neuronales	16	lbfgs	tanh	(11,9,7,5,3)	0.88782224	0.04312936	Redes Neuronales	52	adam	identity	(11,9,7,5,3)	0.86119867	0.05247567
Redes Neuronales	17	lbfgs	tanh	(13,11,9,7,5,3)	0.88298376	0.04840562	Redes Neuronales	53	adam	identity	(13,11,9,7,5,3)	0.8595517	0.0562678
Redes Neuronales	18	lbfgs	tanh	(21,15,9,3)	0.88432028	0.04271349	Redes Neuronales	54	adam	identity	(21,15,9,3)	0.8604346	0.05262875
Redes Neuronales	19	lbfgs	relu	(5,3)	0.76988594	0.16781753	Redes Neuronales	55	adam	logistic	(5,3)	0.88346259	0.04356749
Redes Neuronales	20	lbfgs	relu	(7,5,3)	0.83562671	0.09348695	Redes Neuronales	56	adam	logistic	(7,5,3)	0.88342494	0.03952955
Redes Neuronales	21	lbfgs	relu	(9,7,5,3)	0.87978031	0.04448801	Redes Neuronales	57	adam	logistic	(9,7,5,3)	0.88190435	0.04279432
Redes Neuronales	22	lbfgs	relu	(11,9,7,5,3)	0.8752064	0.0446602	Redes Neuronales	58	adam	logistic	(11,9,7,5,3)	0.88614594	0.04052862
Redes Neuronales	23	lbfgs	relu	(13,11,9,7,5,3)	0.52003181	0.04006363	Redes Neuronales	59	adam	logistic	(13,11,9,7,5,3)	0.88243527	0.04055279
Redes Neuronales	24	lbfgs	relu	(21,15,9,3)	0.88662863	0.04011936	Redes Neuronales	60	adam	logistic	(21,15,9,3)	0.88252373	0.03983591
Redes Neuronales	25	sgd	identity	(5,3)	0.86256137	0.05282942	Redes Neuronales	61	adam	tanh	(5,3)	0.88135704	0.03772752
Redes Neuronales	26	sgd	identity	(7,5,3)	0.85989317	0.05408672	Redes Neuronales	62	adam	tanh	(7,5,3)	0.88024338	0.04232707
Redes Neuronales	27	sgd	identity	(9,7,5,3)	0.85973702	0.05387656	Redes Neuronales	63	adam	tanh	(9,7,5,3)	0.88303656	0.04330461
Redes Neuronales	28	sgd	identity	(11,9,7,5,3)	0.86030911	0.05257009	Redes Neuronales	64	adam	tanh	(11,9,7,5,3)	0.88316276	0.04036142
Redes Neuronales	29	sgd	identity	(13,11,9,7,5,3)	0.86190868	0.05384885	Redes Neuronales	65	adam	tanh	(13,11,9,7,5,3)	0.87531416	0.05213801
Redes Neuronales	30	sgd	identity	(21,15,9,3)	0.86048838	0.05306772	Redes Neuronales	66	adam	tanh	(21,15,9,3)	0.87867699	0.03932291
Redes Neuronales	31	sgd	logistic	(5,3)	0.86124723	0.05138885	Redes Neuronales	67	adam	relu	(5,3)	0.88166351	0.04571206
Redes Neuronales	32	sgd	logistic	(7,5,3)	0.7698591	0.08254627	Redes Neuronales	68	adam	relu	(7,5,3)	0.87380403	0.04242091
Redes Neuronales	33	sgd	logistic	(9,7,5,3)	0.72989652	0.07037349	Redes Neuronales	69	adam	relu	(9,7,5,3)	0.88311551	0.03740613
Redes Neuronales	34	sgd	logistic	(11,9,7,5,3)	0.45320431	0.0683945	Redes Neuronales	70	adam	relu	(11,9,7,5,3)	0.87960204	0.0369097
Redes Neuronales	35	sgd	logistic	(13,11,9,7,5,3)	0.3591786	0.14671174	Redes Neuronales	71	adam	relu	(13,11,9,7,5,3)	0.87237822	0.05291438
Redes Neuronales	36	sgd	logistic	(21,15,9,3)	0.30764215	0.10930984	Redes Neuronales	72	adam	relu	(21,15,9,3)	0.88415551	0.03962361

Apéndice 5B – Extracción de Variables relevantes

```
## RANDOM FOREST
# 2.MEJOR MODELO

from sklearn import model_selection
from sklearn.ensemble import RandomForestClassifier

criterion="gini"
n_estimator=60
max_depth=10
min_samples_split=200
clf = sk.ensemble.RandomForestClassifier(criterion=criterion, n_estimators=n_estimators, max_depth=max_depth, min_samples_split=min_samples_split)
scores = sk.model_selection.cross_val_score(clf, X, y, scoring="roc_auc", cv=5)
print(scores.mean(), scores.std(), '\n')

import numpy as np
clf.fit(X, y)
importances=clf.feature_importances_
std = np.std([tree.feature_importances_ for tree in clf.estimators_],
              axis=0)

indices = np.argsort(importances)[::-1]

# Print the feature ranking
print("Feature ranking:")

for f in range(X.shape[1]):
    print("%d. feature %d (%f)" % (f + 1, indices[f], importances[indices[f]]))
```

Referencias

- ⁱ Patrutiu-Baltes, L. (2016). Inbound Marketing-the most important digital marketing strategy. *Bulletin of the Transilvania University of Brasov. Economic Sciences. Series V*, 9(2), 61.
- ⁱⁱ Mitik, M., Korkmaz, O., Karagoz, P., Toroslu, I. H., & Yucel, F. (2017). *Data Mining Approach for Direct Marketing of Banking Products with Profit/Cost Analysis. The Review of Socionetwork Strategies*, 11(1), 17–31. doi:10.1007/s12626-017-0002-5
- ⁱⁱⁱ Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). The KDD process for extracting useful knowledge from volumes of data. *Communications of the ACM*, 39(11), 27-34.
- ^{iv} LAMB, Charles W.; HAIR, Joe F.; MCDANIEL, Carl. *Essentials of marketing*. Cengage Learning, 2011.
- ^v Piatetsky-Shapiro, G. (1991). Discovery, analysis, and presentation of strong rules. *Knowledge discovery in databases*, 229-238.
- ^{vi} Aggarwal, C. C. (2015). *Data mining: the textbook*. Springer
- ^{vii} Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2016). *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann.
- ^{viii} Weiss, S. I., and Kulikowski, C. 1991. *Computer Systems That Learn: Classification and Prediction Methods from Statistics, Neural Networks, Machine Learning, and Expert Systems*. San Francisco, Calif.: Morgan Kaufmann.
- ^{ix} Jain, A. K., and Dubes, R. C. 1988. *Algorithms for Clustering Data*. Englewood Cliffs, N.J.: PrenticeHall.
- ^x Agrawal, R.; Mannila, H.; Srikant, R.; Toivonen, H.; and Verkamo, I. 1996. Fast Discovery of Association Rules. In *Advances in Knowledge Discovery and Data Mining*, eds. U. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, 307–328. Menlo Park, Calif.: AAAI Press.
- ^{xi} Rodrigo C. Barros, André C.P.L.F de Carvalho, Alex A. Freitas, “Automatic Design of Decision-Tree Induction Algorithms”, 2015 - ISBN 978-3-319-14230-2, 978-3-319-14231-9
- ^{xii} Sahu, S. K., Prasad, M. B. N. V., & Tripathy, B. K. A Support Vector Machine Binary Classification and Image Segmentation of Remote Sensing Data of Chilika Lagloon.
- ^{xiii} Reyzin, L., & Schapire, R. E. (2006, June). How boosting the margin can also boost classifier complexity. In *Proceedings of the 23rd international conference on Machine learning* (pp. 753-760). ACM.
- ^{xiv} Fawcett, T. (2006). An introduction to ROC analysis. *Pattern recognition letters*, 27(8), 861-874.
- ^{xv} <https://www.kaggle.com/contact>
- ^{xvi} STEINMANN, Michael; MATEI, Sorin Adam; COLLMANN, Jeff. A theoretical framework for ethical reflection in big data research. En *Ethical Reasoning in Big Data*. Springer, Cham, 2016. p. 11-27.