



# Analista de Big data

**Competencia 2:** Segmentación de clientes de acuerdo a su comportamiento como consumidor.

Índice:

1. *Análisis de técnicas de minería de datos.*
2. *Aplicación de minería de datos en diferentes sectores.*
  - *Big data en el sector Retail.*
  - *Big data en el sector de Seguros*
  - *Big data en el sector de la banca*
  - *Big data en Telecomunicaciones*
  - *Big data en el sector transporte.*
  - *Big data en el sector de medios de comunicación*
  - *Big data en el sector de la investigación científica*
  - *Big data en el sector de las ciencias sociales*



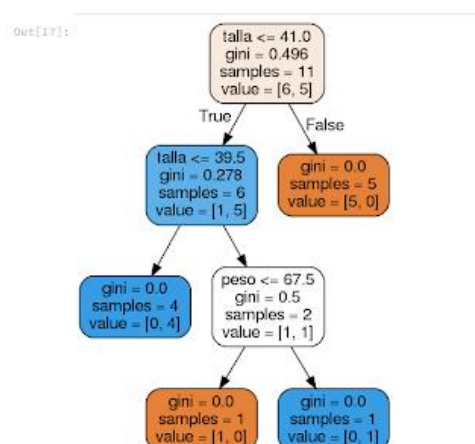
## Análisis de técnicas de minería de datos.

La minería de datos se define como una serie de técnicas encaminadas a identificar patrones implícitos dentro de grandes conjuntos de datos, con el fin de entender sus mecanismos de comportamiento, su interrelación y su potencial evolución futura. El data mining es una tecnología exploratoria clave en los proyectos de Big Data, y se puede poner en funcionamiento tanto para resolver preguntas específicas como para la extracción de información de manera general, buscando tendencias y anomalías en la muestra.

Si bien ya tuvimos una primera introducción a lo referido al procesamiento de datos masivos, en este apartado vamos a analizar un poco mas desde el lado de los diferentes procesos y técnicas que se utilizan dentro del Big data. Vamos a ver 7 técnicas que son utilizadas dentro de este ambiente, teniendo en cuenta que cada una de ellas se ajustara mejor a nuestro caso de aplicación o no, es decir que las técnicas aquí propuestas tendrán un mejor desempeño según el problema que estemos abordando.

### 1. Árboles de decisión.

Los árboles de decisión son diagramas lógicos que plantean, ante una determinada situación, cuáles son las opciones de intervención posibles, agregando sus implicaciones, costes, ventajas y desventajas. Se basan en la aplicación de un algoritmo clasificador que, a partir de un nodo, desarrolla ramas (decisiones) y determina el potencial resultado de cada una de ellas.

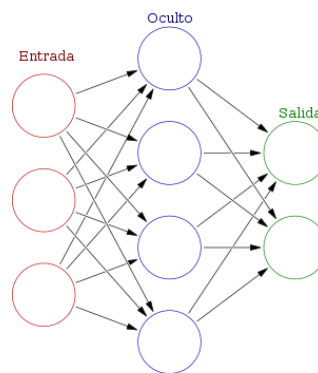


Ejemplo de árbol decisión



## 2. Redes neuronales.

Las redes neuronales son modelos que, a través del aprendizaje automático, intentan llenar los vacíos de interpretación en un sistema. Para ello imitan, en cierto modo, las conexiones entre neuronas que se producen en el sistema nervioso de los seres vivos. Las redes neuronales se engloban dentro de las técnicas predictivas de minería de datos y, como todo modelo de machine learning, es preciso entrenarlas con distintos data sets con los que ir matizando los pesos de las neuronas para asegurar la fiabilidad de sus respuestas. Existen diferentes tipos de redes neuronales para data mining, como el perceptrón simple y la multicapa.



Red neuronal simple

## 3. Clustering.

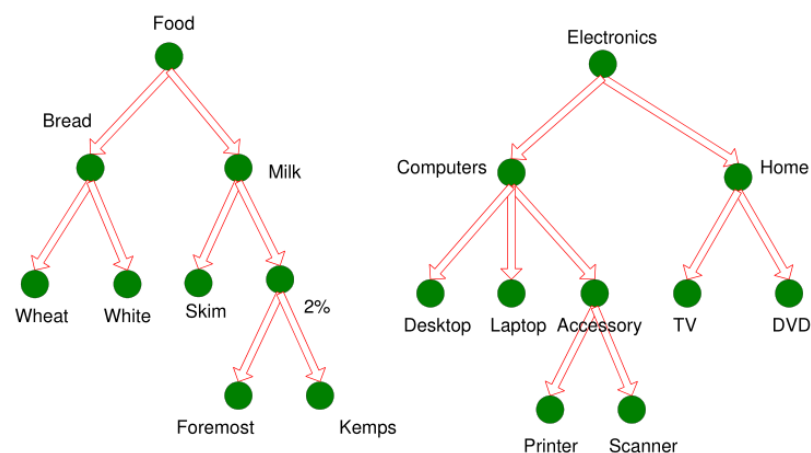
El clustering o agrupamiento en minería de datos tiene como objetivo la segmentación de elementos que presentan alguna característica definitoria en común. En este caso, el algoritmo atiende a condiciones de cercanía o similitud para hacer su trabajo. Esta técnica de data mining está muy extendida en el mundo del marketing para el envío de correos y promociones personalizadas a los usuarios que integran una base de datos.



Clustering o agrupamiento

#### 4. Extracción de reglas de asociación.

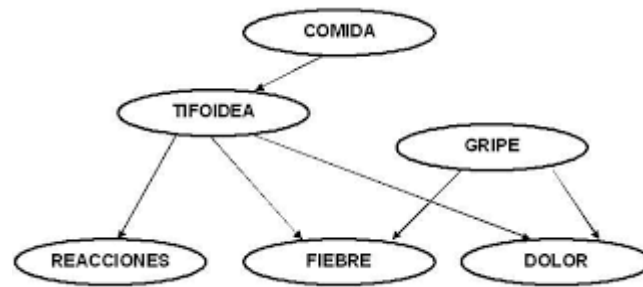
La extracción de reglas de asociación como técnica de minería de datos persigue la inferencia de silogismos del tipo si.../entonces... a partir de conjuntos de registros. Esta búsqueda de regularidades nos permite discriminar y conocer mejor a una muestra, y establecer qué atributo o combinación de atributos es probable que traiga consigo una determinada consecuencia.



Reglas de asociación

#### 5. Redes bayesianas.

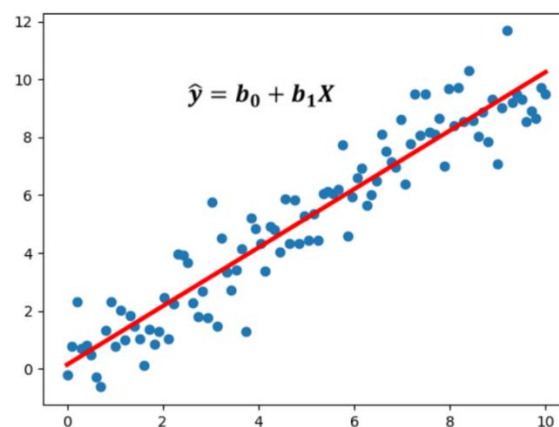
Las redes bayesianas son representaciones gráficas de relaciones de dependencia probabilística entre distintas variables. Sirven para dar solución tanto a problemas descriptivos como predictivos. Entre sus aplicaciones se incluyen el diagnóstico médico o el cálculo del riesgo en el sector financiero y asegurador.



Redes bayesianas

## 6. Regresión.

La regresión como técnica de minería de datos toma como punto de partida una serie histórica para, a partir de ella, predecir qué sucederá a continuación. De manera resumida, podemos decir que, a través de este método, se localizan regularidades dentro de los datos que permiten trazar una línea de evolución extrapolable al futuro.



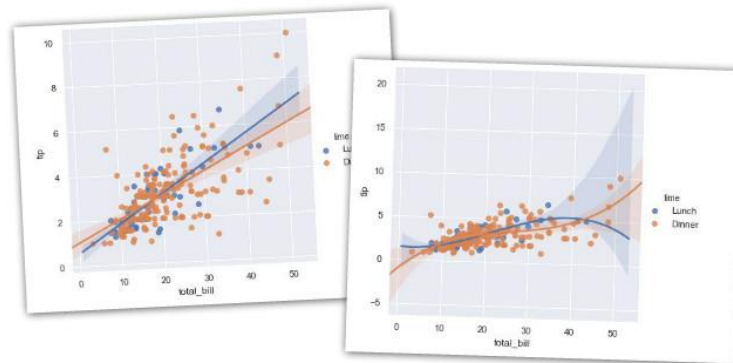
Regresión (lineal)

## 7. Modelo estadístico.

El modelado estadístico pretende dibujar el mapa de relaciones entre variables explicativas y dependientes, y mostrar cómo cambia a medida que lo hacen los parámetros considerados en su planteamiento. Lo hace estableciendo una ecuación matemática que intenta reproducir la realidad de la manera más fiel



posible, incorporando, incluso, la influencia del azar y el posible margen de error en el cálculo.



Modelos estadísticos

## Aplicación de minería de datos en diferentes sectores.

El tratamiento de grandes volúmenes de información o Big Data tiene como objetivo resolver problemas de negocio. Vamos a ver varios casos de usos reales donde se utilizaron técnicas de Big Data e información estructurada y no estructurada en distintos sectores y vamos a focalizarnos en su impacto en el negocio, Cabe destacar que si bien no son casos reales, son de total aplicación en cada sector.

- Big data en el sector Retail.

Una multinacional mayorista quería conocer mejor su canal de distribución, las tiendas donde vendían sus productos. Entonces comenzaron a trabajar con la información interna que tenían, qué productos vendían, en qué tiendas e intentar entender cuáles eran sus mejores canales de distribución y cuáles no lo eran tanto. Pero, esto no es suficiente porque el problema de utilizar solo la información interna es que se conoce qué pasa en el resto de tiendas en las que no tienes productos o qué pasaría en las tiendas si colocaras otros productos distintos a los que actualmente vendes.



Para hacer todo esto se juntó información estructurada y no estructurada, interna y externa para conseguir tener una visión global del entorno de cada una de las tiendas. Se juntaron datos del catastro, del padrón, de Google maps, de movilidad, de tránsito en la zona y con esto se consiguió tener una imagen clara y definida de la zona alrededor de cada una de las tiendas (el entorno de cada tienda). Si a esto lo unimos a la información que tenemos internamente sobre las ventas de las tiendas entonces somos capaces de conocer para las tiendas propias todos sus productos vendidos, toda su información y todo su potencial. Incluso el potencial de esas tiendas si distribuyésemos otro tipo de productos para poder utilizar la cesta de productos que distribuíamos en esa tienda. Y para aquellas tiendas que actualmente no distribuían productos lo que se podía saber es cuál eran los productos que mejor resultado iban a tener por el entorno que tenían y además cuál era el potencial de cada tienda.

Para hacer todo este proyecto se involucran a muchísimas áreas de negocios, de estrategia, insight, comercial, logística. Las técnicas que se utilizaron no fue solo una, se utilizaron técnicas no supervisadas de clustering para poder agrupar tiendas similares, con características similares y también se utilizaron técnicas de supervisión o regresión para poder calcular todo el potencial de cada una tienda. Además, se utilizaron técnicas de inferencia para poder establecer valores desconocidos en tiendas de las que teníamos poca información. El impacto del negocio fue muy alto. El crecimiento medio de las tiendas era del dos por ciento. Pues, cuando se aplicaron sobre un conjunto de tiendas las recomendaciones que salían de los análisis realizados se pudo crecer 10 puntos por encima del resto y esto traducido a dinero es muchísima cantidad de dinero. Lo importante fue no solo extraer el conocimiento utilizando fuentes internas y fuentes externas sino poner en valor el modelo y hacer acciones concretas sobre tiendas concretas. Además, también se creó una metodología para seleccionar las nuevas tiendas de distribución para poder comercializar los productos y servicios de la multinacional mayorista desde un punto de vista óptimo.



## Caso de uso en el sector Retail

<b>Caso</b>	Segmentación de cadena de tiendas	<b>Áreas involucradas</b>	❖ Estrategia ❖ Insights ❖ Comercial ❖ Logística
<b>Objetivo</b>	Optimizar la red de distribución y detectar nuevos puntos de distribución eficientes	<b>Técnicas utilizadas</b>	Aplicación de modelos algorítmicos para la predicción de ventas potenciales y futuras mediante regresión e inferencia
<b>Descripción</b>	Clusterización de los establecimientos en función de indicadores clave del público objetivo en base a su comportamiento	<b>Fuentes de datos</b>	<b>Datos internos:</b> base de datos de ventas. <b>Datos externos:</b> movilidad (SmartSteps), AEMET, padrón, Catastro, Open Street Maps...
<b>Desafíos</b>	Implementar las mejores estrategias y tácticas en cada tienda sobre el surtido, precio, material y acciones comerciales. Para ello, es imprescindible una integración geoespacial complete de fuentes de información y ponerlo en valor mediante un framework de visualización que permita ver representados los resultados de forma sencilla para una toma de decisiones ágil y eficiente		
<b>Puesta en valor</b>	En los puntos de venta en los que se ha aplicado el modelo, se ha crecido 10 puntos por encima sobre el resto cuando el canal muestra un decrecimiento de 2 o 3 puntos de media. Identificación de las tiendas con alto potencial de venta en las que la multinacional mayorista todavía no distribuía el producto		

- Big data en el sector de Seguros.

Una aseguradora tenía un problema con el fraude. Notaban que tenían más fraude de la media del mercado, y lo que querían era minimizarlo de alguna forma. Para ello, lo que habían detectado algunas veces es que aparecían operaciones de fraude entre personas vinculadas. A veces entre familiares, a veces entre amigos, entre conocidos, o incluso, entre empresas de grupo. El objetivo del análisis era comprender mejor estas relaciones entre los distintos asegurados para evitar el fraude. Claro, tenemos que entender que parte de las personas involucradas en un siniestro eran aseguradoras de la compañía, y otra parte los intervinientes, podían no ser clientes, por lo que la información de esas personas no estaba disponible internamente. Lo que se hizo fue incluir la información interna de los asegurados y después se completó con información externa, catastro, páginas amarillas, registro mercantil, el objetivo era realizar una red de relaciones que nos permitiese conocer los patrones de comportamiento y detectar esas relaciones que estaban relacionados con los siniestros fraudulentos. Para ello, se utilizaron técnicas de teoría de grafos sobre bases de datos de grafos, para dibujar y trabajar el mapa de relaciones. Además, se incluyeron modelos gráficos probabilísticos para el análisis de la influencia, para intentar identificar cuáles eran las personas que eran más influyentes con respecto a los demás. Por último, se realizaron modelos supervisados para calcular ese “scoring” de relación entre distintos intervinientes, con el objetivo de descartar las casualidades y focalizarse en





esos fraudes donde había personas relacionadas de algún tipo. A veces eran familiares directos, a veces eran familiares indirectos, a veces eran amigos, y otras veces simplemente es que estaban o pertenecían a la misma asociación. Ahí el uso de redes sociales y las relaciones entre ellos, ayudó muchísimo a detectar todas estas relaciones.

Entender el mapa de relaciones de las personas y la influencia de los factores sociodemográficos ha sido una de las claves y el éxito de este proyecto. Además, lo que nos dimos cuenta es que realmente había que hacer dos análisis independientes. Uno, por un lado, en las personas físicas, y otro por otro lado, las personas jurídicas. Porque los orígenes y las fuentes de información eran diferentes. Igual que las personas físicas se utilizaban más redes sociales, en las personas jurídicas se utilizaban más datos de registros mercantiles o empresas de datos financieros y fiscales de las compañías. Obviamente, las personas jurídicas hay personas que trabajan como administradores y como empleados, que después tenemos que incluir las relaciones, por lo que se hizo un metamodelo por encima que incluía todos los datos anteriores. El impacto del negocio fue brutal. Se detectaron muchísimas relaciones que no existían, que permitieron detectar muchísimo el fraude, y el fraude se redujo significativamente. Además, se creó un ranking de los clientes con más influencia para poder orientar todas las acciones de la compañía hacia esos “influencers” para que ellos pudiesen, con su influencia, impregnar a los demás.

## Caso de uso en el sector Seguros

<b>Caso</b>	<b>Detección del fraude</b>	<b>Áreas involucradas</b>	❖ Riesgos ❖ Altas ❖ Sinistros
<b>Objetivo</b>	Identificación de relaciones ocultas entre titulares y tomadores para evitar estafas en la contratación	<b>Técnicas utilizadas</b>	Teoría de Grafos para analizar el mapa de relaciones, modelos gráficos probabilísticos para el Análisis de la influencia y modelos supervisados para el Scorings de relación
<b>Descripción</b>	Cruce masivo de base de datos para la detección de núcleos familiares, no familiares y empresariales en función de diversas variables: dirección, apellidos, datos de contacto, números de pólizas, datos empresariales, etc	<b>Fuentes de datos</b>	<b>Datos internos:</b> base de datos de clientes <b>Datos externos:</b> INE, Páginas Blancas, Catastro, Páginas Amarillas, BBDD Financieras, Registros de deuda...
<b>Desafíos</b>	Entender el mapa de relaciones de las personas y la influencia de los factores sociodemográficos en estas relaciones utilizando teorías de redes y grafos. Necesidad de realizar estudios independientes para personas físicas y para personas jurídicas y combinarlos para una completitud máxima		
<b>Puesta en valor</b>	<b>Detectadas relaciones no declaradas entre más del 16% de los registros de asegurados analizados</b> <b>Ranking de los clientes que más influencia poseen en el mercado de la aseguradora</b>		



- Big data en el sector de la banca.

Una entidad financiera quería detectar sus clientes de alto valor para poder focalizar todos sus esfuerzos comerciales sobre ellos. Para ello, contaba con muchísima información interna de transacciones, histórico del TPV, CRM, etcétera, y un histórico muy amplio, de forma que fue fácil mediante técnicas descriptivas detectar y determinar cuáles son los clientes de alto valor. Pero ahora surgían dos problemáticas. Uno, ¿cómo calculamos el potencial de aquellos que no son clientes? Y dos, de aquellos que son clientes, pero que sus posiciones son muy reducidas, ¿cómo podemos determinar el valor? Pues claro, si un cliente pertenece a otra entidad, no puedes identificarlas. Entonces, aquí surgieron una serie de técnicas de machine Learning para poder identificar clientes parecidos para poder conocer el comportamiento de los mismos, con la premisa de que, si los clientes son parecidos, tendrán un potencial parecido. Pero esto no era suficiente. Entonces, lo que se hizo fue incluir información externa tanto de catastro como redes sociales, como de otras fuentes de datos públicas, que nos permitiesen enriquecer la visión que teníamos de las personas y de las relaciones, porque si una persona está relacionada con clientes de alto valor, probablemente también será de alto valor. Con este objetivo se estuvo trabajando en toda la gesta de información y en toda la fusión de integración de fuente y después en la construcción de esos modelos de Machine Learning. Por un lado, no supervisados para clasificar y segmentar a los clientes, y, por otro lado, supervisado para estimar ese valor. Al final, se consiguió mejorar el éxito de la selección de los clientes de alto valor en un 12,3% consiguiendo mejorar la cartera de clientes premium mucho más de lo que se estaba mejorando en los meses previos. De hecho, se consiguieron más de 189.000 clientes nuevos de alto valor con posiciones significativas en la entidad.



## Caso de uso en el sector Banca

<b>Caso</b>	<b>Detección de clientes de valor</b>	<b>Áreas involucradas</b>	❖ CRM ❖ Estrategia ❖ Comercial
<b>Objetivo</b>	Identificación de clientes potenciales de alto valor no identificados por la entidad financiera	<b>Técnicas utilizadas</b>	Análisis con técnicas de Machine Learning no supervisado e inferencia de variables en redes de Markov
<b>Descripción</b>	Análisis avanzado de información interna. Enriquecimiento con el valor de la vivienda e inmuebles a partir de datos externos, relaciones de redes sociales y datos públicos para inferir variables económicas y financieras de no clientes	<b>Fuentes de datos</b>	<b>Datos internos:</b> transacciones, histórico TPVs, CRM <b>Datos externos:</b> Catastro, redes sociales, Boletín Oficial del Estado
<b>Desafíos</b>	Realizar estimaciones de variables financieras y económicas de personas de las que no se disponía de ninguna información financier mediante la inferencia de estos valores de los clientes a los que se les integraba información interna estructurada y externa no estructurada. Entre estas fuentes, evaluar el uso de red social profesional como potencial proxy de valor para cualificar a los clientes en base a su puesto, empresa, formación, etc.		
<b>Puesta en valor</b>	<b>Aumento de un 12,3% de éxito de selección</b> de cartera de clientes Premium <b>Más de 189.000 nuevos clientes con valoración financieras muy elevadas</b>		

- Big data en Telecomunicaciones

Uno de los problemas que más preocupa en el sector Telecomunicaciones es el churn o el abandono de clientes, es decir, aquellos clientes que cambian de compañía. El problema que tiene el churn, es que es de golpe sin que se pueda apreciar un cambio en el comportamiento del consumo del cliente porque lo único que hace es cambiar de una compañía a otra. Para poder analizar este problema, se consideró toda la información interna que se tenía. Las altas, las bajas, los pagos, las transacciones para ver si había cambios en el comportamiento de los clientes que nos pudiese ayudar a predecir ese churn. Pero lo interesante no fue solamente predecir quiénes serían las personas que iban a abandonar para poder actuar sobre ellas, sino lo interesante en este caso fue entender cuáles fueron los motivadores que hacían que los clientes abandonasen la compañía, ¿para qué? Para poder cambiar internamente los procesos de la compañía para mejorar ese proceso y por lo tanto la experiencia del cliente y así mitigar el abandono de los clientes. Las técnicas que se utilizaron fueron técnicas de algoritmos supervisados y no supervisados sobre todo con un enfoque estadístico, para poder entender las relaciones entre las variables y además también se hicieron árboles de decisión con el objetivo de conocer las sendas de desvinculación para poder atajar el problema en las fases tempranas. Para hacer este proyecto se involucraron áreas como CRM, marketing, recuperación, operaciones, etcétera y lo que se pretendía era mejorar la experiencia del usuario y así mitigar el churn. Y esto se conseguía entendiendo realmente los motivadores del abandono mediante estas técnicas. Al final una vez puestas las recomendaciones en producción de los cambios



que había que hacer, lo que se midió fue el éxito precisamente de estos cambios y se cuantificaron 3.8 millones de euros de beneficios en dos años, gracias a las recomendaciones que salieron del conocimiento extraído de estos modelos.

## Casos de uso en Telecomunicaciones

<b>Caso</b>	<b>Predicción del churn</b>	<b>Áreas involucradas</b>	❖ CRM ❖ Marketing ❖ Admisión ❖ Seguimiento
<b>Objetivo</b>	Conocer los motivadores del abandono de los clientes para optimizar los procesos de concesión y seguimiento	<b>Técnicas utilizadas</b>	Algoritmos supervisados y no supervisados para el churn score. Árboles y Bosques para las sendas de desvinculación y motivadores de impago.
<b>Descripción</b>	Clustering de los clientes con mayor probabilidad de fuga según diferentes hábitos de consumo. Análisis de las transiciones entre los clusters para detectar los motivadores de fuga de clients y las sendas de desvinculación de estos para actuar en su mitigación	<b>Fuentes de datos</b>	<b>Datos internos:</b> altas y bajas, atención al cliente, codificación llamadas, incidencias, reclamaciones, saldo... <b>Datos externos:</b> Sociodemográficos
<b>Desafíos</b>	Analizar el grafo de relaciones de los clientes cuyo tamaño y complejidad requirió una plataforma tecnológica de gran tamaño. Operativizar los modelos construidos y que se explotasen en un tiempo reducido que permitiese detectar el abandon en una fase temprana para poder gestionar el abandon con suficiente antelación		
<b>Puesta en valor</b>	<b>Incremento acumulado de ingresos estimado en 3,8 millones de euros en dos años</b> tras la implantación del modelo. <b>Mejora en el segmento objetivo</b> sobre el que aplicar las diferentes acciones de retención e identificación de <b>patrones de comportamiento</b> para optimizar los planes de retención.		

- Big data en el sector transporte.

Una aerolínea quería conocer mejor a sus clientes para poder hacerle ofertas personalizadas para que mejorase la experiencia de usuario; ofrecerles vuelos a determinados sitios, ofrecerles complementos, determinado tipo de asiento, si mejor en ventanilla o pasillo, en función a sus gustos y preferencias, para simplificar así todo el proceso del viaje. Para ello, se contaba con toda la información interna de los vuelos previos, las evaluaciones, las encuestas, las compras, las navegaciones por la web, etcétera, de todos estos viajeros. Se involucró a distintas áreas, el área de CRM, el área de Marketing digital, y muy especialmente al área de Experiencia de usuario, porque el objetivo era facilitar todo el proceso a los usuarios. Mediante técnicas de Big Data y Machine Learning, lo que siguió, lo que se consiguió fue hacer una microsegmentación muy personalizada de los gustos y preferencias de los usuarios para poder personalizar todos los momentos del viaje, desde la compra, el embarque, el asiento, el vuelo y toda la experiencia resultante previa. En impacto fue el aumento de un 237% del número de clientes que contrataban productos a través de sus campañas personalizadas y la generación de nuevas variables





para más de 840.000 clientes, porque lo que se hizo fue inferir variables en función a las variables que teníamos de nuevos clientes. Con todo esto, lo que se consiguió fue aumentar la ratio de satisfacción de los clientes y la valoración que los clientes hacían de la aerolínea.

## Caso de uso en el sector Transporte

<b>Caso</b>	<b>Segmentación avanzada de clientes</b>	<b>Áreas involucradas</b>	❖ CRM ❖ Marketing Digital ❖ Experiencia de Usuario
<b>Objetivo</b>	Caracterización y segmentación de clientes para disponer de una visión 360°, aumentar la retención y satisfacción de los mismos.	<b>Técnicas utilizadas</b>	Técnicas de Machine Learning no supervisadas combinadas con técnicas de regresión y diseño de experimentos para realizar análisis causa efecto
<b>Descripción</b>	Clusterización de la base de datos y posterior segmentación para el lanzamiento de campañas a partir de nuevas variables: emails, origen del dato, Facebook, LinkedIn, histórico de vuelos e información del canal de venta con el objetivo de mejorar la experiencia tanto durante el vuelo como en los periodos previos y posteriores.	<b>Fuentes de datos</b>	<b>Datos internos:</b> campañas de captación, web, encuestas, compras,...
<b>Desafíos</b>	Se presentaron dos grandes desafíos, por un lado, mejorar la experiencia, mediante recomendaciones personalizadas, de los usuarios puntuales de los que no se disponía prácticamente de información y por otro lado, diferenciar entre el uso profesional y el uso personal en la personalización de la experiencia del usuario		
<b>Puesta en valor</b>	<b>Aumento de 237% de campañas personalizadas para la gestión de clientes</b> <b>Generación de nuevas variables</b> para más de 840K clientes y potenciales Hasta 200K nuevos <b>clientes personalizados (+55%)</b>		

- Big data en el sector de medios de comunicación

El impacto del Big Data en el mundo de la comunicación audiovisual es muy importante, por un lado, por la capacidad de manejar grandes cantidades de datos, y por el otro por la cantidad de información disponible hoy en día. Es importante porque las personas que acceden a los contenidos audiovisuales, como ser emisiones de televisión, son millones de personas que acceden a los contenidos. Con lo cual tenemos un montón de impactos de datos a guardar, de cosas a manejar, que sirven para conocer los gustos de la audiencia, poder aconsejar, poder seguirlos.

Manejar todos estos datos, extraer características, sacar métricas, poder elaborar recomendaciones, es muy difícil. Otro apartado es la gran capacidad de proceso, que sirve también para poder analizar imágenes de vídeo y poder analizar contenidos de audio. De manera que podamos hacer clasificaciones automáticas y trabajos automáticos. Por tanto, los dos ámbitos, la capacidad para manejar contenidos y la capacidad para procesar, son relevantes e



importantes para el sector de las comunicaciones. Se requieren profesionales que tengan unas características determinadas. Uno, es que tengan una formación en telecomunicaciones, en informática, en matemática o en física. Que tengan también interés en estas materias. No es un tipo de profesional que esté normalmente en el mercado porque las escuelas forman ingenieros, o forman físicos, o forman matemáticos, y aquí requiere una cierta especialización. El proceso en si es “Metadatear” (que sería como la acción de clasificar y etiquetar los contenidos) todo esto, en el sentido de que, si vas a buscar alguna cosa el archivo, has de saber qué es lo que hay ahí. Por tanto, hay que reconocer las imágenes de las personas que están allí, los textos. “Metadatear” todo esto de una manera automática es un proyecto importante para poder avanzar hacia la automatización en ese sector. Otro proyecto importante podría ser buscar la relación uno a uno entre el emisor y la audiencia. Llevar una relación uno a uno entre el emisor y la audiencia, quiere decir tener varios millones de clientes a los que has de atender cada día. Por tanto, estos dos proyectos nada más ya son dos elementos que necesitan estos pilares de Big Data. La relación uno a uno con la audiencia y el poder sacar metadatos, y poder indexar todo el material que se va archivando. Hay algunos temas más que van de la base de nuevas generaciones de procesadores que traen inteligencia y que son capaces, por ejemplo, en un partido de fútbol de conocer los planos y las secuencias, y hacerte propuestas en función del audio, de editados, de resúmenes. O hay tecnologías que trabajan a partir de bases de datos estructurados, por ponerte, textos y gráficos.

- Big data en el sector de la investigación científica.

Quería hablaron un poco sobre la relación entre Big data y el mundo de la investigación científica. Las experimentaciones que se hacen en física, astrofísica, en biología, en astronomía, en observación de la tierra, generan grandes cantidades de datos. Esas cantidades de datos empezaron a acumularse, empujando los límites de la tecnología de computadores en los años 70 en el mundo, por ejemplo, de las físicas de partículas en laboratorios de aceleradores como puede ser el CERN. A lo largo de las décadas se han ido perfeccionando estos experimentos y acumulan cada vez más y más datos. De hecho, en este momento el Gran colisionador de hadrones, en el CERN en Ginebra ya se ha acumulado más de 200 petabytes de datos, que se



corresponden a miles de millones de colisiones subatómicas. Evidentemente estos datos hay que analizarlos y hay que sacar señales que cada vez son más difíciles de buscar porque están enterradas dentro de procesos físicos que conocemos y entonces hemos de buscar lo que no conocemos.

Otro ejemplo puede ser la astronomía que hoy en día se mezcla con la astrofísica y con la cosmología. Queremos saber básicamente de dónde viene nuestro universo. Y para hacer eso, lo que hacemos es observar con telescopios que hoy en día son robotizados y están todas las noches de todo el año observando el universo y generando también petabytes de datos. Una cosa que hace solo 10 años se pensaba imposible pero hoy en día es cotidiana. Un experimento es el “Dark Energy Survey”, que ya ha acumulado miles de millones de objetos medidos del universo y de hecho abre una nueva frontera que es el hacer estudios estadísticos sobre lo que vemos en el universo. Estos datos no son tan grandes como los de física de partículas, pero son extremadamente complejos y lo que queremos es buscar correlaciones entre ellos, correlaciones que no son tan diferentes de alguna de las correlaciones que tenemos en otros ámbitos totalmente diferentes como pueden ser las redes sociales, otro ámbito en el cual se utiliza el Big data.

Otro ejemplo es el de la biología y lo que queremos estudiar es estudiar la estructura de la vida. Para ello contamos con máquinas que secuencian los genomas y si en el pasado, la primera secuenciación del genoma se tardó casi 10 años, hoy en día podemos secuenciar un genoma en un día o menos. Por tanto, estamos acumulando datos, una vez más, extremadamente complejos.

Habíamos hablado de que las personas debían “conocer el negocio” o tener un conocimiento en el ámbito de la tarea que desarrolla, para este caso lo que se busca son personas que tengan cierta transversalidad con respecto a digamos los pilares clásicos de la ciencia. En particular, buscamos gente que tenga bastantes conocimientos matemáticos porque contrario de lo que se pueda pensar, en el Big data, tenemos mucha modelización matemática. No es solo meter los datos en un ordenador y ya está. Sino que necesitamos modelizar estos datos para que sea posible entonces extraer información útil y veraz a una velocidad razonable de estos datos. Entonces, matemáticas es uno de los pilares.

- Big data en el sector de las ciencias sociales.



Parece evidente que, en el tema del gran cambio que implica internet y el mundo digital, uno de los factores que cada vez está adquiriendo más peso son las perspectivas que se abren en relación a la gran cantidad de información que se va acumulando debido, tanto a las propias búsquedas de la gente en los buscadores como en la propia movilidad de la gente, simplemente, utilizando las aplicaciones que hay de mapas y de movilidad, simplemente, en las actividades de consumo. Es decir, la capacidad de almacenar datos que deriven de las actuaciones y de las actividades de la gente en su quehacer ordinario, que antes para poder determinar pautas de conducta y de actuación de las personas tenías que acudir, o bien, a encuestas o a elementos de carácter estadístico que no siempre estaban disponibles, ahora, con la capacidad de almacenamiento de datos y de uso de estos datos que nos dan las redes digitales, es evidente que se está abriendo un campo de análisis, de experimentación, de trabajo, de utilización enorme. Es evidente, por ejemplo, que algunas compañías de distribución de agua, o energéticas, están cambiando los contadores tradicionales, por ejemplo, por contadores digitales. Esto, aparentemente, es útil porque no necesitan llevar a alguien a inspeccionar y analizar los consumos de cada persona en ese sitio determinado y, por tanto, se ahorran, hay un elemento de eficiencia que es el ahorro que implica que no tengas que acudir a ese sitio. Pero, al margen de esto, la gran importancia que tienen estos contadores digitales es que permiten saber, por ejemplo, para una compañía de aguas, de distribución de aguas, cuáles son los consumos de agua en toda la ciudad por ramas de actividad, por horarios de utilización del agua, por momentos del año. Y esto permite una planificación, tanto de la propia distribución del agua como de la venta de esos datos a posibles usuarios que quieran, por ejemplo, instalar una panadería. Si tú quieres instalar una panadería en una ciudad, probablemente, la agencia proveedora de agua es la que tendrá más información sobre cuál es el mejor sitio para que tú vayas a situar tu panadería. ¿Por qué? Porque, a través del control que hagan de los contadores digitales, sabrán quién está en el fondo vendiendo más pan porque está utilizando más agua, y a qué horas lo hace, y qué días de la semana, etcétera, etcétera. Si tú utilizas el sistema de aplicaciones móviles para comprar comida por las noches y que te la traigan a tu casa, aparentemente, las compañías que están utilizando, que están prestando este servicio, están, simplemente, moviendo comida del sitio que la producen al sitio que la consumen. Pero, la cantidad de datos que están acumulando sobre el tipo de comida que cada día en el mundo se está pidiendo, seguramente, les hace los mejores proveedores de información para que alguien





que quiera abrir un restaurante sepa si lo que se está llevando ahora es el tártaro de salmón, el ceviche o la pizza cuatro quesos. Porque el nivel de información al día, al momento, en el segundo en que se está produciendo, la están teniendo ellos más que otros ámbitos. Lo mismo ocurriría con la cantidad de datos que nosotros, simplemente, moviéndonos con nuestro smartphone, estamos generando. Esto puede permitir, por ejemplo, que haya estudios que relacionen la movilidad de las personas con la salud de esas personas. Se habla mucho de los efectos terapéuticos que tiene andar, se utilizan bases de datos de carácter estadístico y estudios que son de un altísimo coste en muchos casos. Realmente, la utilización de esta información acumulada por parte de las compañías de telefonía, tendría una utilización clara en el ámbito de la salud. Es decir que estamos, simplemente, en los inicios de lo que puede ser una revolución o, de lo que es, ya una revolución de carácter productivo, laboral, social que tiene, evidentemente, aspectos que pueden ser muy positivos y aspectos que no lo pueden ser tanto. ¿Quién tiene más información hoy sobre el nivel de movilidad en las ciudades, las compañías de transporte público o Google Maps y Waze? Seguramente, Google Maps y Waze porque tienen control constante, a través del sistema de análisis de GPS, de la movilidad de las personas. Es una frase que se dice muchas veces cuando tú utilizas Google Maps, estás utilizando un instrumento gratuito, esto quiere decir que tú eres la mercancía. En este caso tú eres la mercancía, ¿por qué?, porque estás incorporando información que alguien puede utilizarla.