

Big Data, IA y Machine Learning



red.es



"El FSE invierte en tu futuro"
Fondo Social Europeo

Índice

1. Big Data

- 1.1 Computer Science
- 1.2 Redefiniendo Big Data

2. Data Science: AI y ML

- 2.1 AI y ML. ¿son lo mismo?
- 2.2 ¿Puede un ordenador pensar?
- 2.3 ML: Carácter predictivo
- 2.4 Industrias de aplicación ML
- 2.5 La importancia de los algoritmos
- 2.6 Proceso Machine Learning
- 2.7 Aprendizaje supervisado vs no
- 2.8 ELT y Data Lakes
- 2.9 Ejemplos y reflexiones
- 2.10 ¿Qué resolvemos con ML?
- 2.11 Clasificación
- 2.12 Regresión
- 2.13 Clustering

3. ML: Entender el proceso

- 3.1 Pregunta adecuada
- 3.2 Identificación y preparado de datos
- 3.3 Identificar y aplicar algoritmos
 - 3.3.1 Árboles
 - 3.3.2 Bayes
 - 3.3.3 Redes neuronales
 - 3.3.4 kNN
- 3.4 Evaluación y ajuste de modelos
- 3.5 Uso y presentación de modelo

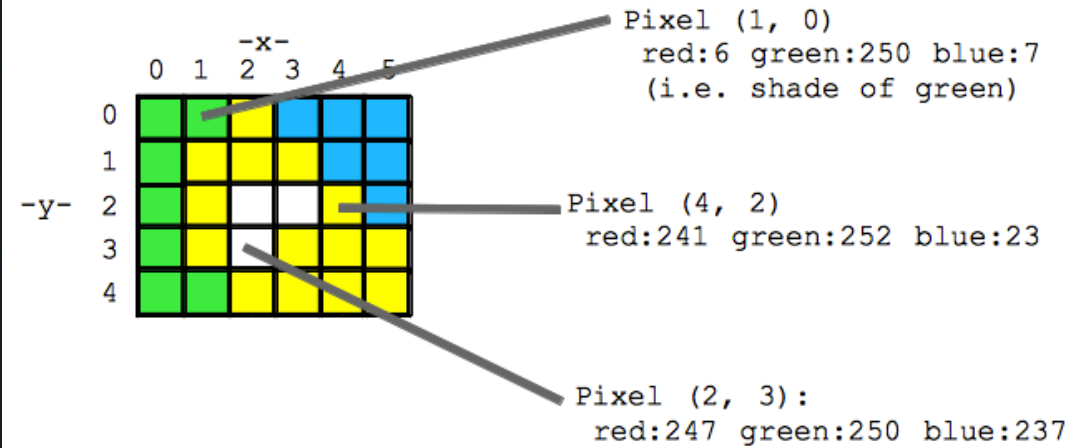
4. Ejercicios

5. Examen

1. Big Data

BIG DATA, IA Y MACHINE LEARNING

R0	R0	R55	R48	R39	R17	R11	R31	R0	R0
G0	G0	G160	G182	G183	G119	G68	G43	G0	G0
B0	B0	B146	B135	B119	B128	B134	B113	B0	B0
R0	R79	R58	R53	R86	R24	R1	R3	R23	R0
G0	G132	G205	G240	G246	G42	G20	G69	G99	G0
B0	B188	B142	B124	B142	B219	B199	B121	B66	B0
R77	R8	R133	R166	R145	R67	R7	R0	R4	R33
G51	G231	G250	G208	G130	G57	G10	G90	G133	G71
B216	B42	B178	B220	B249	B244	B228	B114	B39	B55
R68	R48	R164	R200	R167	R76	R9	R0	R0	R19
G45	G185	G218	G192	G152	G65	G13	G91	G148	G89
B228	B127	B211	B250	B251	B245	B231	B115	B21	B43
R45	R76	R8	R166	R141	R65	R7	R0	R0	R15
G32	G118	G231	G156	G126	G55	G10	G90	G148	G96
B226	B210	B42	B248	B249	B244	B227	B114	B21	B39
R12	R22	R61	R133	R64	R15	R0	R0	R6	R29
G13	G22	G148	G250	G85	G16	G5	G55	G83	G64
B209	B237	B174	B178	B220	B235	B209	B126	B75	B81
R10	R1	R7	R8	R6	R33	R0	R1	R20	R41
G14	G6	G120	G231	G192	G156	G10	G6	G47	G67
B178	B212	B133	B42	B68	B131	B176	B150	B117	B85
R25	R3	R0	R0	R0	R0	R1	R14	R40	R55
G29	G6	G54	G179	G190	G144	G27	G37	G82	G64
B133	B175	B146	B36	B27	B48	B131	B122	B105	B70
R0	R19	R4	R1	R0	R6	R20	R40	R56	R0
G0	G21	G9	G104	G115	G29	G47	G83	G77	G0
B0	B132	B146	B58	B46	B122	B118	B107	B83	B0
R0	R0	R40	R19	R15	R26	R42	R58	R0	R0
G0	G0	G40	G34	G20	G43	G66	G69	G0	G0
B0	B0	B94	B102	B112	B107	B94	B75	B0	B0



1.1 Big Data y Computer Science

Antes de hablar de conceptos como big data, machine learning, inteligencia artificial o ciencia de datos, es imprescindible conocer el término de ‘computer science’ o ‘ciencias de la computación’.

En este sentido, según la Wikipedia:

*“Las **ciencias de la computación** o **ciencias de la informática** son las ciencias formales que abarcan las bases teóricas de la información y la computación, así como su aplicación en sistemas computacionales. El cuerpo de conocimiento de las ciencias de la computación es frecuentemente descrito como el estudio sistemático de los procesos algorítmicos que describen y transforman información: su teoría, análisis, diseño, eficiencia, implementación, algoritmos sistematizados y aplicación”*

Así pues, como veremos durante las próximas diapositivas, el Big data no es más que un campo dentro de todo lo que conocemos como ciencias de la computación o computer science

1.2 Redefiniendo Big Data

Ya hemos hablado de la importancia de 3 elementos básicos que conforman el Big Data:

- Velocidad a la que se consume la información
- Variedad de información
- Volumen de información

Pero todo esto, no nos sirve de nada si no incorporamos un componente de:

VALOR

Si los datos no sirven para aprender, descubrir o analizar, todo lo anterior no sirve de nada.

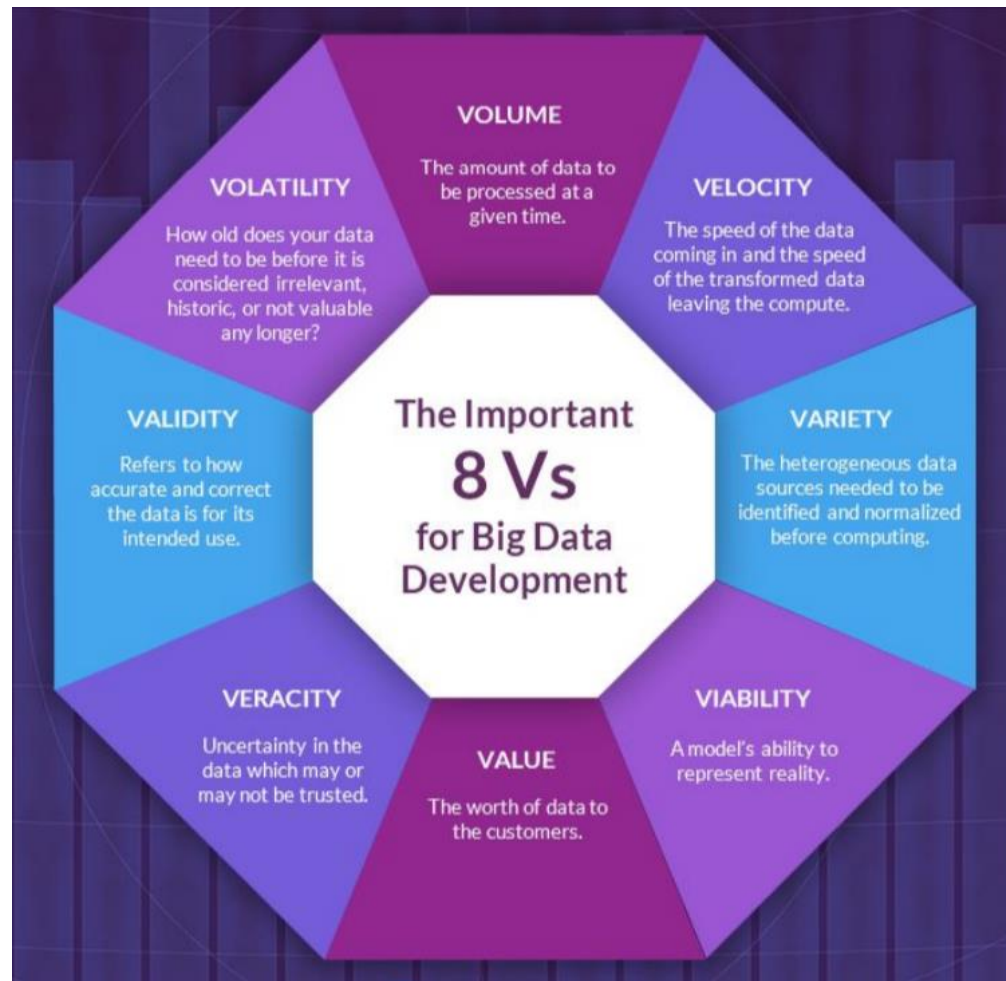
Recoger información es importante pero por sí solos, los datos no nos ayudarán en nada. Es importante no caer en el *hype* del marketing y la innovación asociado al término de Big Data. El big data es algo mucho más amplio y profundo que solamente extraer y almacenar datos.

1.2 Redefiniendo Big Data

¿4 V's u 8 V's?

- ✓ Volumen
- ✓ Variedad
- ✓ Velocidad
- ✓ Valor

- ✓ Viabilidad
- ✓ Veracidad
- ✓ Validez
- ✓ Volatilidad



2. Data Science: AI y Machine Learning

2.1 AI y ML: ¿Son lo mismo?

Técnicamente diríamos que el machine learning y la inteligencia artificial son prácticamente lo mismo. Realmente, conviven en distintos niveles uno dentro del otro. A grandes rasgos podríamos concluir que:

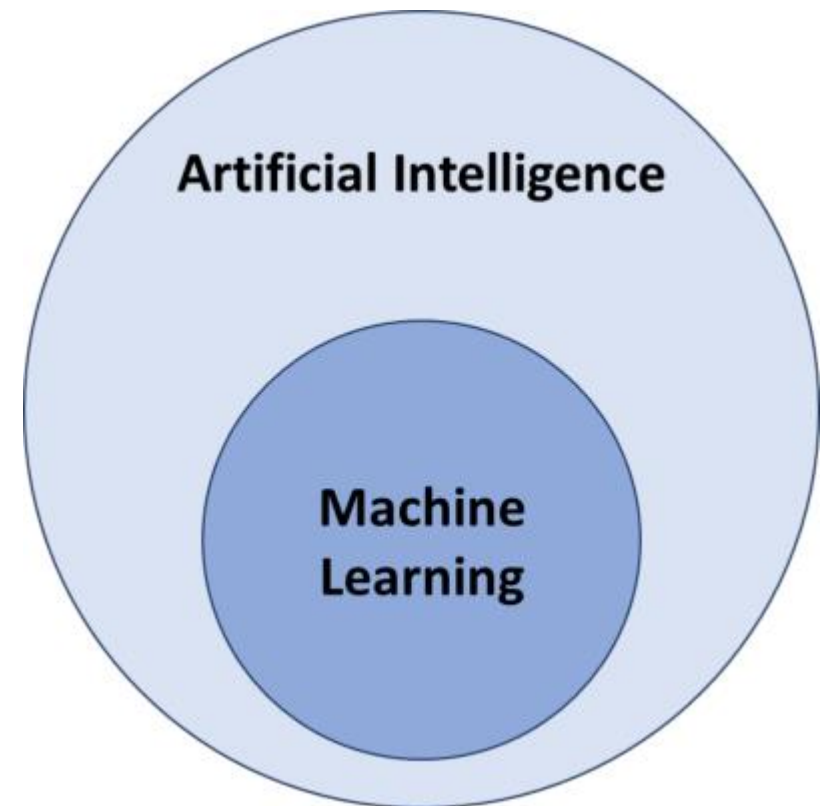
Inteligencia artificial

Entendemos como IA el hecho de que las máquinas imiten comportamientos parecidos al de un ser humano. Puede abarcar muchísimas formas y aplicaciones.

Machine learning

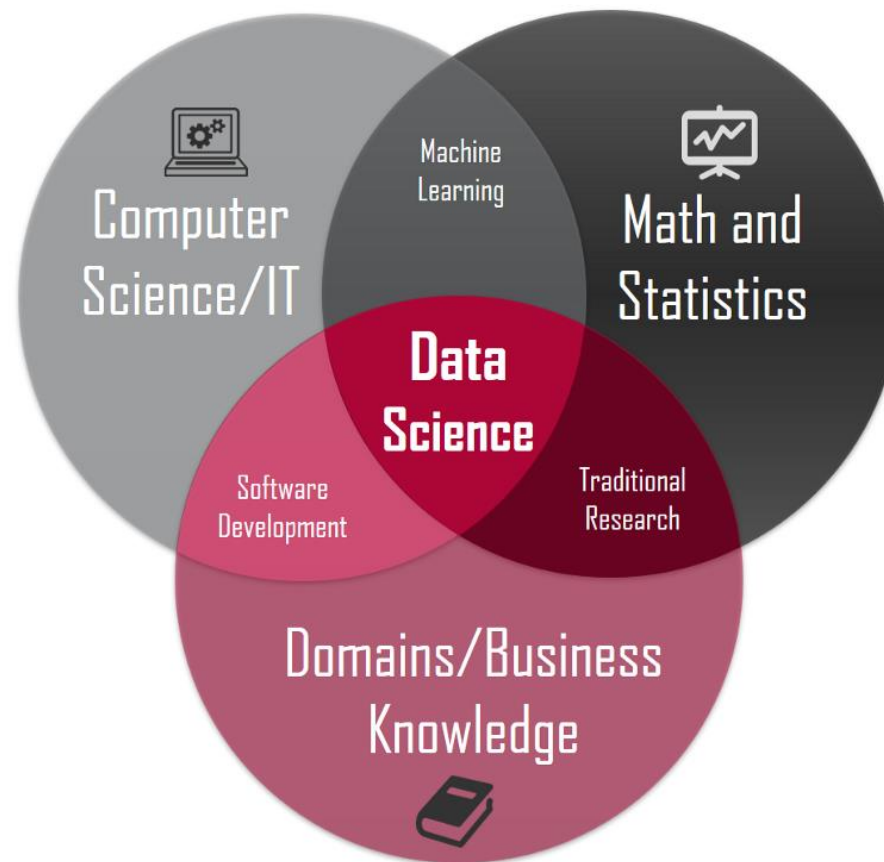
Lo entendemos como un subconjunto de inteligencia artificial en el cual las máquinas aprenden y se entrenan para comportarse como humanos y hacer predicciones.

Deep learning



2.1 AI y ML: ¿Son lo mismo?

Entonces...¿cómo se relacionan estos términos con las ciencias de la computación, el software, la matemática, la estadística o la ciencia de datos? Fijémonos en la siguiente imagen (*cotejar con lo visto en la sesión 1*):



2.2 ¿Puede un ordenador pensar?

A menudo escuchamos o nos hacemos preguntas de este tipo cuando nos relacionamos u oímos términos como inteligencia artificial. Y a esta pregunta podemos añadir varias reflexiones:

UN ORDENADOR **NO** PUEDE PENSAR COMO UN HUMANO.

UN ORDENADOR **SÍ** PUEDE IMITAR EL COMPORTAMIENTO DE UN HUMANO

En este sentido, la pregunta clave es:

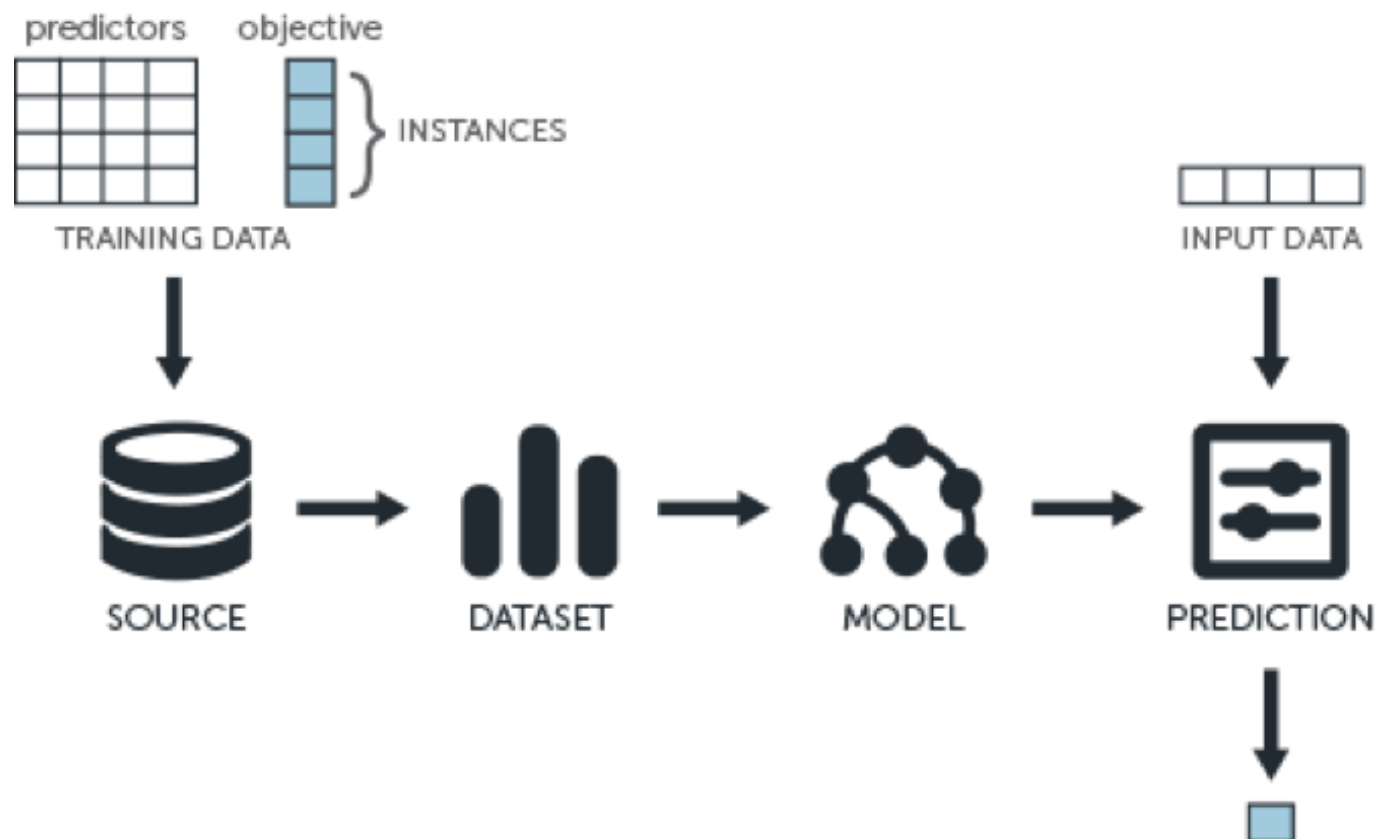
¿REALMENTE IMPORTA SI UN ORDENADOR PUEDA PENSAR COMO UN HUMANO?

- ¿Puede un submarino bucear tal como lo hace un humano o un pez?
- ¿Puede un avión volar tal como lo hace una golondrina?
- ¿Puede un transatlántico nadar como un deportista profesional?

LO QUE IMPORTA NO ES SI PUEDE PENSAR. LO IMPORTANTE ES SI LO PUEDE CONSEGUIR

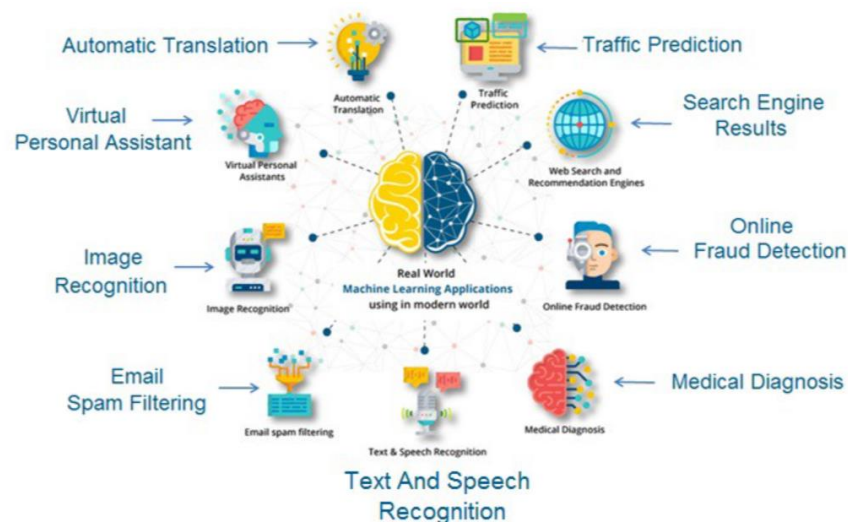
2.3 ML: Carácter predictivo

Cuando hablamos de machine learning, lo que realmente aporta valor e importancia a su existencia y desarrollo no es más que su capacidad para poder aplicar **modelos predictivos**. Esa es la auténtica revolución de esta disciplina, enmarcada dentro de la amplia inteligencia artificial.



2.4 Aplicaciones ML en el día a día

- Coches autónomos
- Reconocimiento facial
- Motores de recomendación
- Detección de fraude en tarjetas de crédito
- Abandono de clientes
- Filtros de spam en los emails
- Categorización y segmentación
- Análisis de emociones
- Personalización de experiencias de clientes
- Predicciones meteorológicas
- Diagnósticos médicos
- Reconocimiento de voz

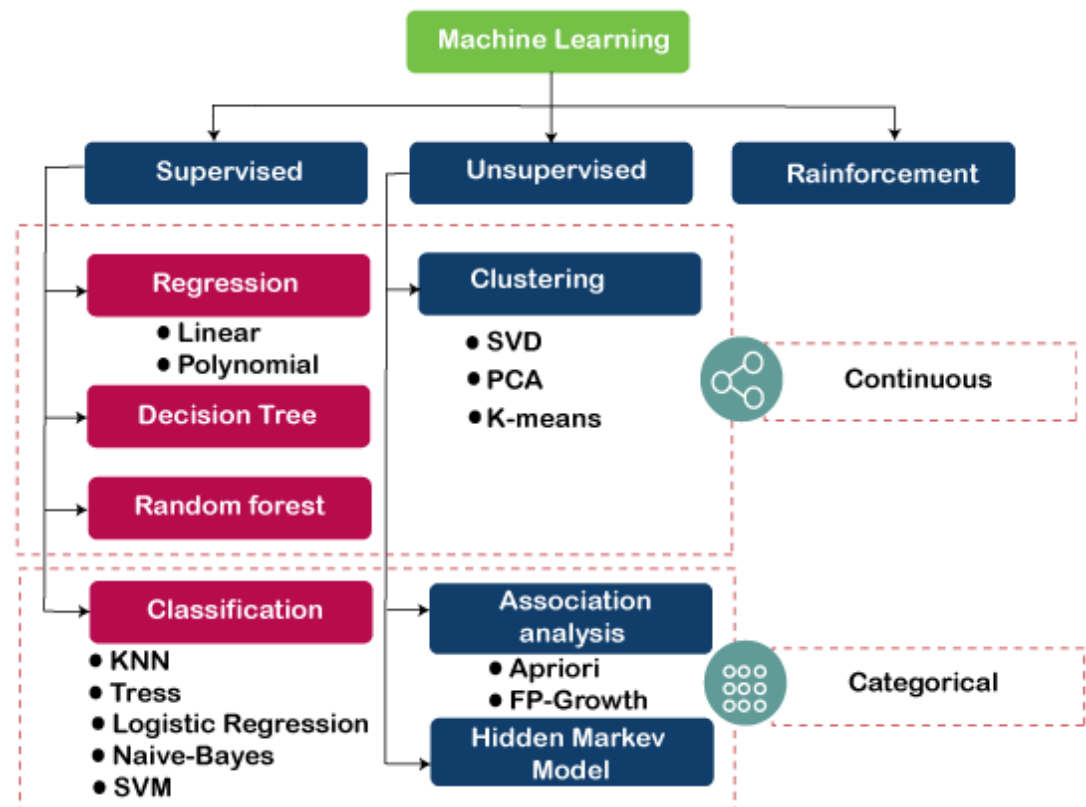


2.5 La importancia de los algoritmos

El machine learning basa prácticamente todo su potencia en la aplicación de algoritmos para aplicar el aprendizaje automático y desarrollar así modelos predictivos.

En este sentido, es sumamente importante la participación de perfiles con conocimientos matemáticos y estadísticos para diseñar estos modelos de aprendizaje y convertirlos en usables y accesibles para la industria.

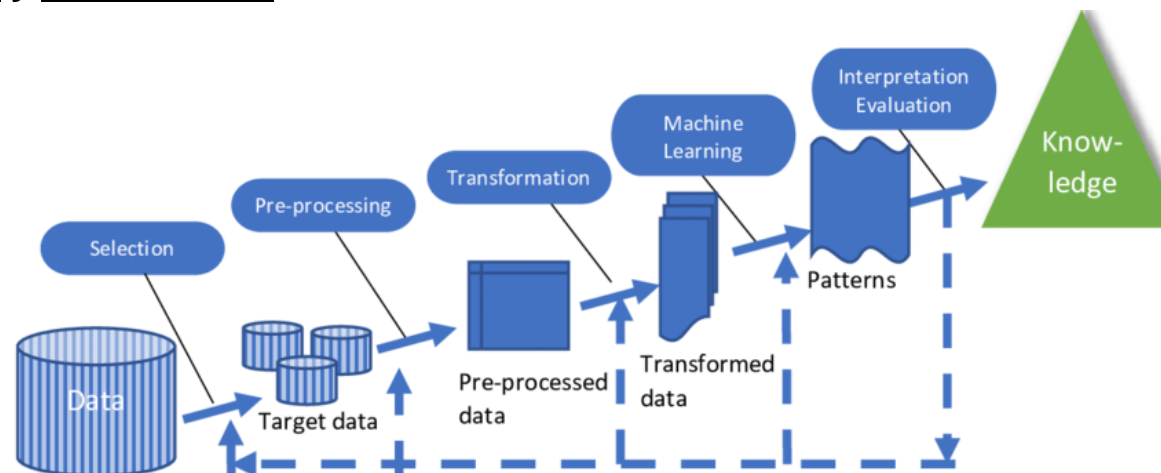
Estos algoritmos pueden estar basados en aprendizaje supervisado o aprendizaje no supervisado como veremos en las siguientes diapositivas.



2.6 Proceso de Machine Learning

Aunque este proceso lo veremos con mayor profundidad en el siguiente capítulo del temario, podríamos decir que todo proyecto de machine learning está basado en la aplicación de estas fases o pasos.

1. Responder a la pregunta adecuada
2. Identificar y preparar los datos correctamente
3. Estudiar y aplicar la mejor solución algorítmica
4. Evaluar y ajustar el modelo para su correcta eficiencia
5. Uso de la solución y presentación del modelo



2.7 Aprendizaje supervisado vs no supervisado

El aprendizaje supervisado es uno de los tipos o modelos más utilizados dentro del machine learning. Para ello, necesitamos contar con unas entradas y salidas de la información sobre la cual vamos a aplicar el aprendizaje automático.

Normalmente, encontramos dos tipos de aprendizaje supervisado:

Clasificación:

Elige entre una lista de opciones previamente definidas y limitada. Por ejemplo elegir un número entre el 0 y el 9. Se puede utilizar en los filtros de spam de los emails.

Regresión:

El objetivo de este algoritmo es predecir números reales o números con infinitas posibilidades. Por ejemplo, lo podríamos aplicar para predecir el precio de un inmueble.

2.7 Aprendizaje supervisado vs no supervisado

Por su parte, en el aprendizaje NO supervisado dentro del machine learning, los algoritmos se aplican para aprender de datos con elementos no etiquetados **buscando** patrones o relaciones entre ellos. En este caso no necesitarían delimitar el número de entradas y salidas.

Existen dos tipos de algoritmos para Machine Learning no supervisado:

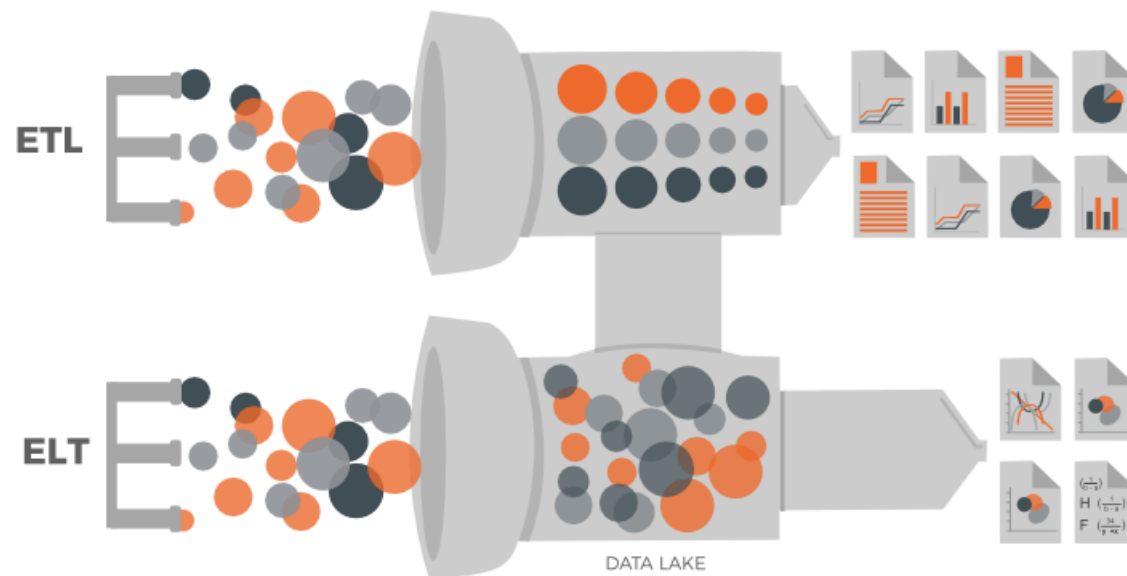
- **Clustering:** clasifica en grupos los datos de salida. Es el caso de las segmentaciones de clientes según qué hayan comprado.
- **Asociación:** descubre reglas dentro del conjunto de datos. Por ejemplo, aquellos clientes que compran un coche también contratan un seguro, por lo que el algoritmo detecta esta regla.

Este tipo de aprendizaje es el que se utiliza por ejemplo en sistemas de recomendación de plataformas de streaming bajo demanda (netflix, hbo, amazon, etc.)

2.8 ELT y Data Lakes

ELT (Extraer, cargar, transformar) es un método diferente de acercarse al flujo de datos, en el que los datos extraídos se cargan primero en el sistema de destino. Las transformaciones se realizan después de que carguemos los datos en el almacén de datos.

Los datos primero se copian en el **data lake** y luego se transforman in situ. Funciona bien cuando el sistema objetivo es lo suficientemente potente como para manejar transformaciones a gran escala. ELT generalmente se usa con bases de datos NOSQL, un dispositivo de datos o una instalación en la nube.



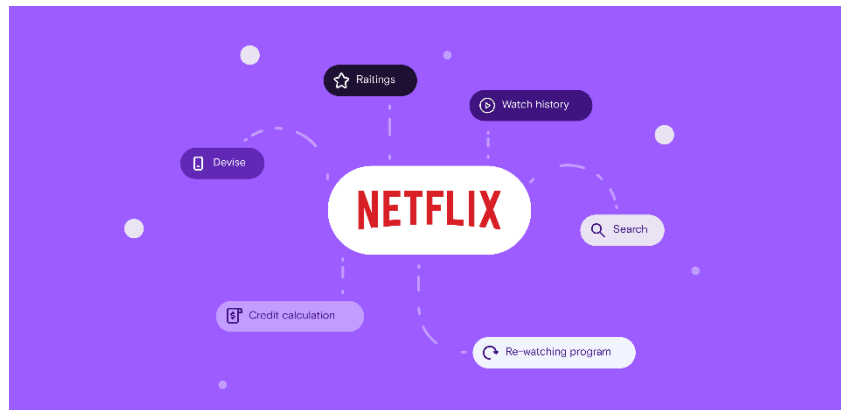
[VER VÍDEO DATADEMIA](#)

2.9 Ejemplos y reflexiones

THISPERSONDOESNOTEXIST

ALGORITMO DE NETFLIX ¿QUÉ TIENE EN CUENTA?

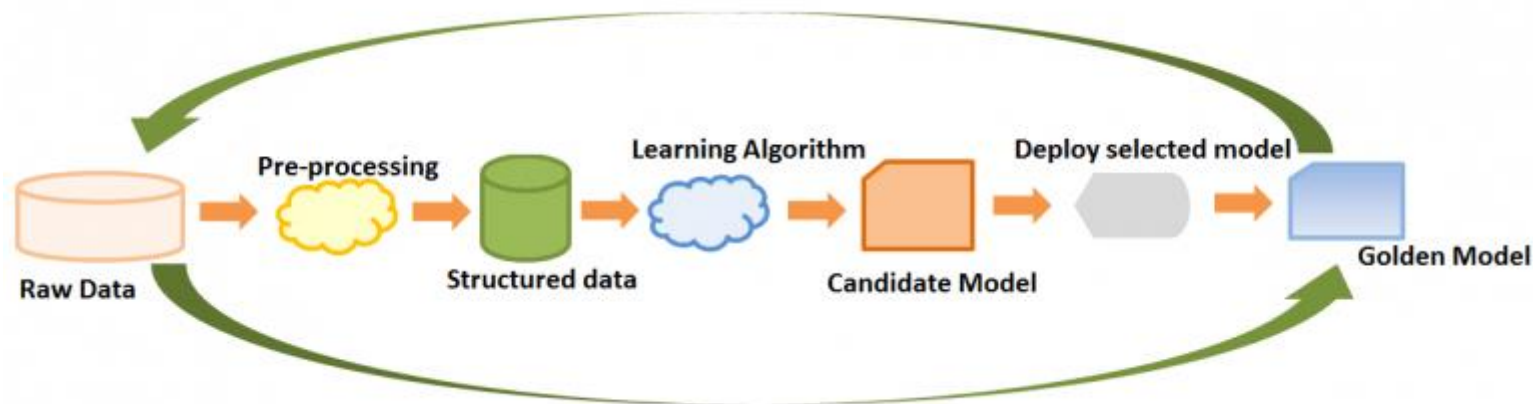
MEDICINA PREDICTIVA



2.10 Qué resolvemos con ML

Llegados a este punto debemos tener en cuenta que el machine learning:

- Se aplica a problemas en los que se ha de aprender a base de datos o ejemplos.
- Se aplican modelos de clasificación, regresión o clustering entre otros.
- No es necesario que los datos sean obligatoriamente estructurados.
- Un elemento clave en la utilización del machine learning es la elección del algoritmo que va a encargarse de aprender y dotar de sentido su aplicación
- Si no tienes los inputs necesarios y las herramientas necesarias no serás capaz de sacar datos relevantes ni resolver el problema que se plantea.

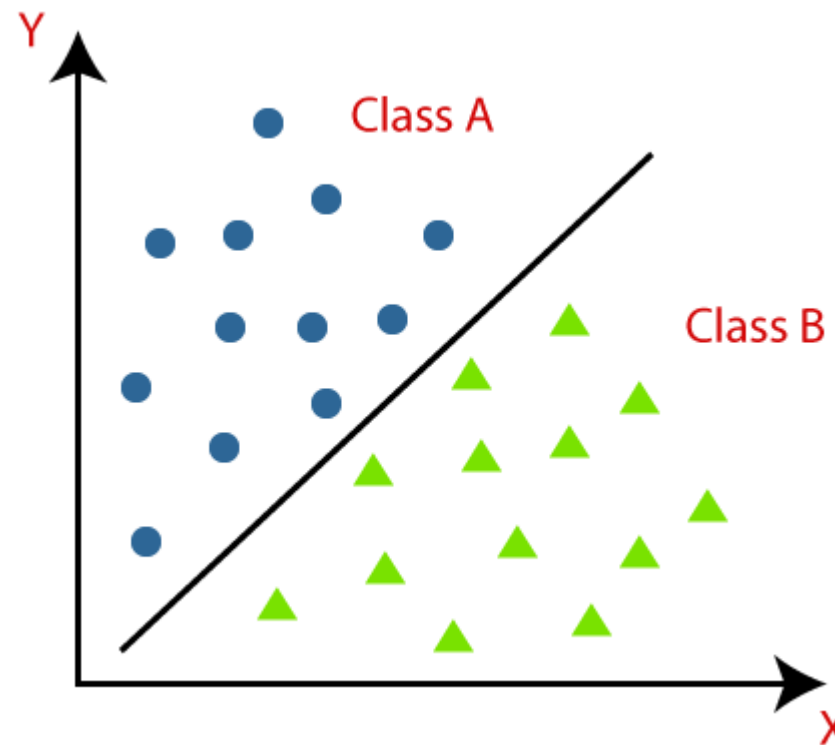


2.11 Clasificación

La clasificación es una subcategoría del aprendizaje supervisado en la que el objetivo es predecir las etiquetas de clase categóricas (discreta, valores no ordenados, pertenencia a grupo) de las nuevas instancias, basándonos en observaciones pasadas.

Clasificación Binaria: Es un tipo de clasificación en el que tan solo se pueden asignar dos clases diferentes (0 o 1). El ejemplo típico es la detección de email spam, en la que cada email es: spam → en cuyo caso será etiquetado con un 1 ; o no lo es → etiquetado con un 0.

Clasificación Multi-clase: Se pueden asignar múltiples categorías a las observaciones. Como el reconocimiento de caracteres de escritura manual de números (en el que las clases van de 0 a 9).



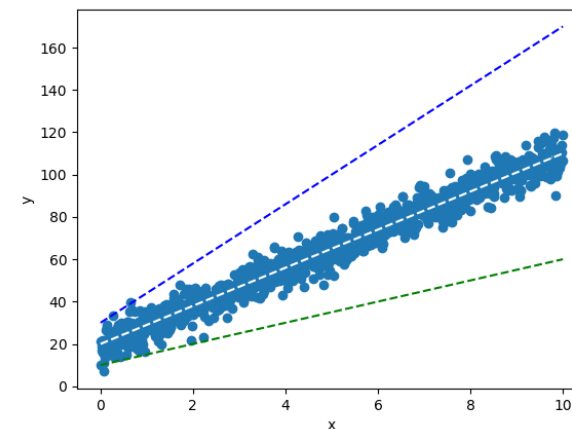
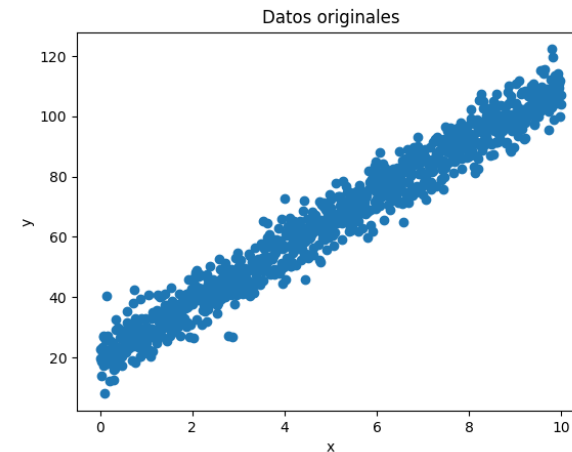
2.12 Regresión lineal

La Regresión Lineal es un procedimiento que permite encontrar la línea recta que mejor representa la relación entre dos variables.

Supongamos que tenemos una serie de datos como la mostrada en la figura, en donde la variable independiente es x y la variable dependiente es y :

Se observa en esta gráfica que la relación entre las dos variables es lineal, es decir que un incremento en la variable x genera un incremento proporcional en la variable y .

El objetivo de la Regresión Lineal es entonces encontrar la línea recta que mejor se ajusta a los datos.



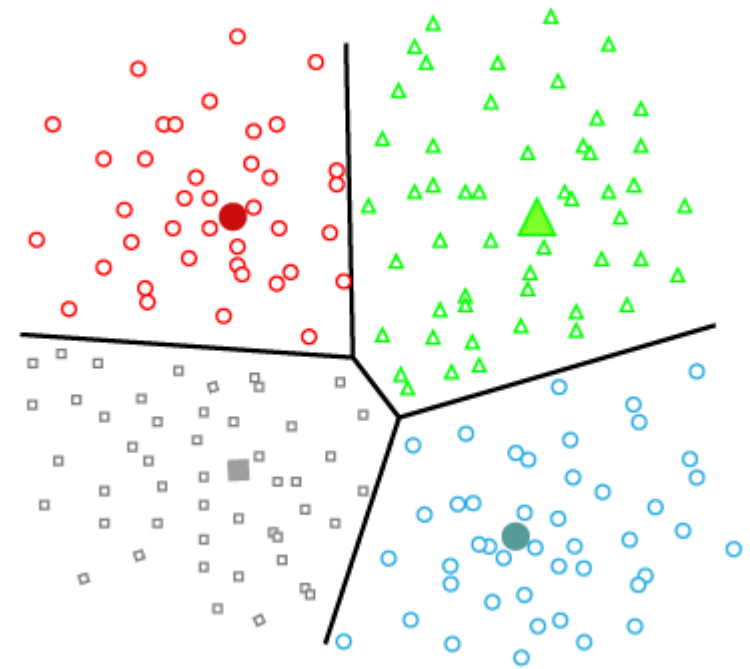
2.13 Clustering

Formalmente, podemos decir que **la agrupación en clústeres** es el proceso de dividir todos los datos en grupos, también conocidos como clústeres, basados en los patrones de los datos.

El Clustering o agrupamiento es una de las formas de aprendizaje no supervisado más utilizada.

Es una gran herramienta para dar sentido a los datos no etiquetados y para agrupar datos en grupos similares.

Un algoritmo de agrupamiento puede descifrar estructuras y patrones en un conjunto de datos que no son aparentes para el ojo humano.



3. ML: Entendiendo el proceso

3.1 La pregunta adecuada

Antes de aplicar cualquier proceso de machine learning o, en general, cualquier proyecto que involucre la extracción, procesamiento y análisis de datos, deberíamos hacernos algunas preguntas que no son en absoluto baladí.

- ¿Qué queremos resolver y qué esperamos obtener de ello?
- ¿Qué recursos tecnológicos tenemos? ¿Contamos con los perfiles y equipo adecuado?
- ¿Cuáles son nuestras fuentes de datos? ¿Son estructuradas o no estructuradas?
- ¿Estamos cumpliendo la legalidad vigente en el uso de estos datos? ¿Nos enfrentamos a información personal sensible? ¿Necesito anonimizar o aplicar una máscara a mis datos?
- ¿Estamos aplicando un sesgo discriminatorio en nuestro enfoque?
- ¿Qué va a definir que el modelo es exitoso o no? ¿Qué porcentaje de error podemos permitirnos?

3.2 Identificar los datos y prepararlos

El segundo paso consiste en conocer los datos o la fuente de datos a utilizar para poder organizarla y prepararla conscientemente para su utilización en el modelo de machine learning.

¿Qué coste tiene acceder a los datos? ¿Vamos a necesitar datos de fuentes externa o son de nuestra propiedad? ¿Qué velocidad de obtención y procesamiento podemos alcanzar?

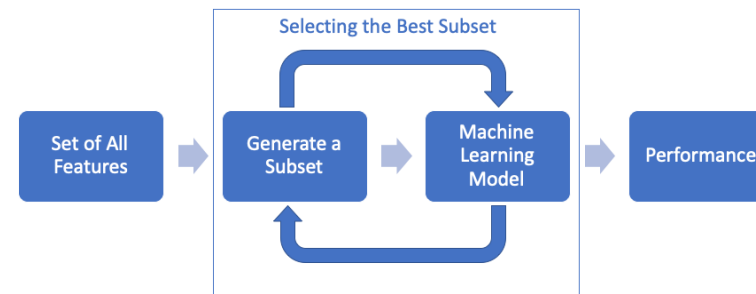
¿Qué variedad de datos vamos a manejar y cuántos formatos vamos a tener que procesar?

A la hora de preparar nuestros datos:

¿Están todos los datos recogidos o hay datos extraviados?

¿Los datos están correctamente estructurados? ¿Hay datos repetidos?

¿Los datos son consistentes?



3.3 Aplicación del algoritmo correcto

El éxito de tu modelo de machine learning va a depender **enormemente** del acierto que se tenga a la hora de relacionar el algoritmo elegido con los datos existentes.

Árboles de decisión

No necesita mucho entrenamiento. Pueden ser de clasificación o de regresión.

Naive Bayes

Es un modelo puramente probabilístico. ¿Qué probabilidad tengo de ver la opción A cuando ocurre B?

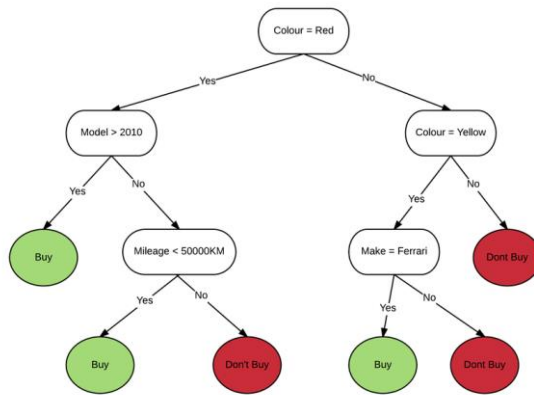
Redes neuronales (VIDEO).

Funcionan como un sistema de neuronas en el cerebro. Todas tienen entradas y salidas y están conectadas entre ellas distribuyendo un “peso” determinado entre cada conexión.

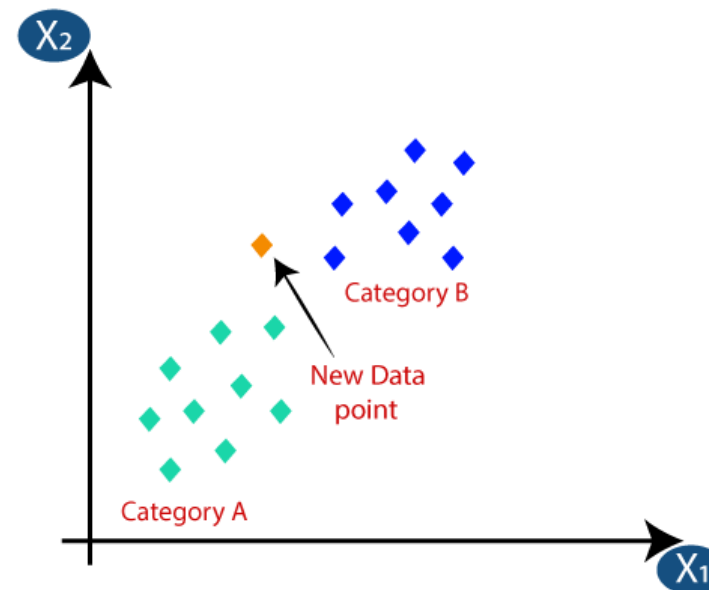
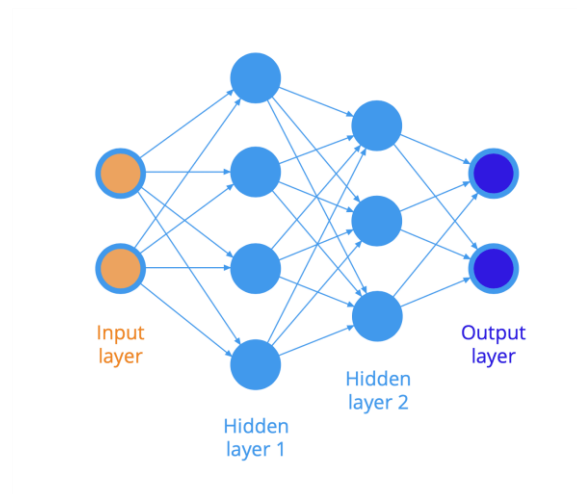
KNN

Sirve esencialmente para clasificar valores buscando los puntos de datos “más similares” (por cercanía) aprendidos en la etapa de entrenamiento

3.3 Aplicación del algoritmo correcto



$$p(C|F_1, \dots, F_n) = \frac{p(C) p(F_1, \dots, F_n|C)}{p(F_1, \dots, F_n)}.$$



3.4 Evaluación y ajuste de modelo

Es el momento de evaluar el rendimiento de tu modelo de machine learning, y en este momento es especialmente crucial la revisión y la corrección del modelo atendiendo a las preguntas que has planteado y respondido en el primer paso del proceso.

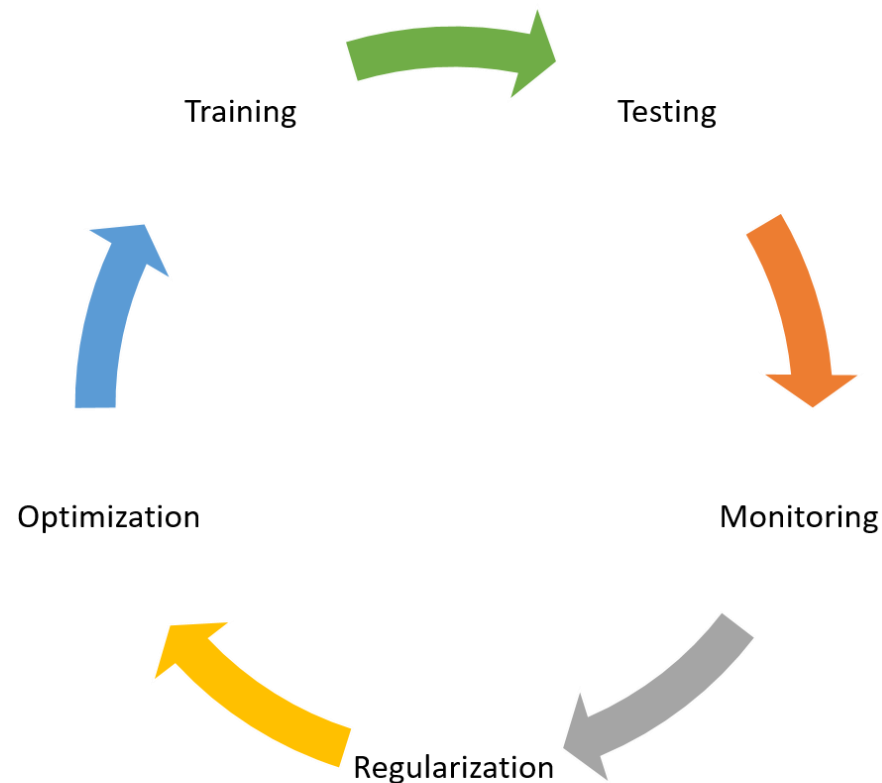
- ¿Necesito ajustar los datos que estoy dando a mi modelo?
- ¿El algoritmo que he elegido satisface y soluciona mi problema?
- ¿Es escalable a la inserción y el análisis de un mayor volumen de datos en el futuro?

- ¿Me está dando resultados muy simples? ¿Puedo percibir la inteligencia aplicada?
- ¿Los resultados son muy complejos? ¿Se aleja demasiado de la realidad y no es comprensible o parecido a un proceso de aprendizaje humano?

SI NO RESPONDES TU PRINCIPAL PREGUNTA O TUS RESULTADOS NO SON LOS ESPERADOS, ES EL MOMENTO DE VOLVER ATRÁS Y MODIFICAR LOS PASOS 2 Y 3.

3.5 Uso y presentación de modelo

Es el momento de empezar a utilizar tu modelo y presentarlo a tu cliente o implementarlo en tu organización. La perspectiva y el juicio humano será la que finalmente confirmará que el modelo funciona y soluciona el problema que se había planteado.



4. Ejercicios

4. Ejercicios

E.1

Visionado de vídeos en los anexos

E.2

Con la ayuda de la extensión 'Data Miner', haz un scraping y descarga los datos de las casas en venta en tu ciudad existentes en la web pisos.com. Ahora, haz una regresión múltiple en una hoja de Excel para predecir un precio en base a m2 y/o número de habitaciones.

E.3

Introducción a Deep Learning y Python: Consulta y explicación del anexo [Google Collab](#) para convertir Celsius a Fahrenheit.

E.4

Separar la clase en dos grupos y aportar al menos 5 argumentos en contra y a favor de usar el big data en el sector de los seguros.

E.5

Examen final



red.es



"El FSE invierte en tu futuro"

Fondo Social Europeo

