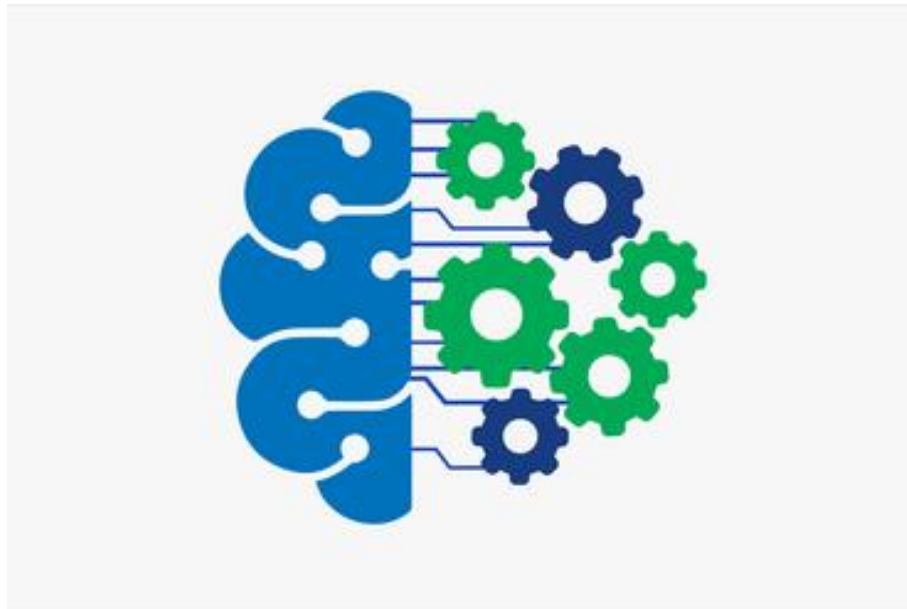


# APLICACIÓN E IMPLEMENTACIÓN DE ALGORITMOS DE MACHINE LEARNING EN LA INDUSTRIA



André Omar Chávez Panduro

Correo : [andrecp38@gmail.com](mailto:andrecp38@gmail.com) / [09140205@unmsm.edu.pe](mailto:09140205@unmsm.edu.pe)

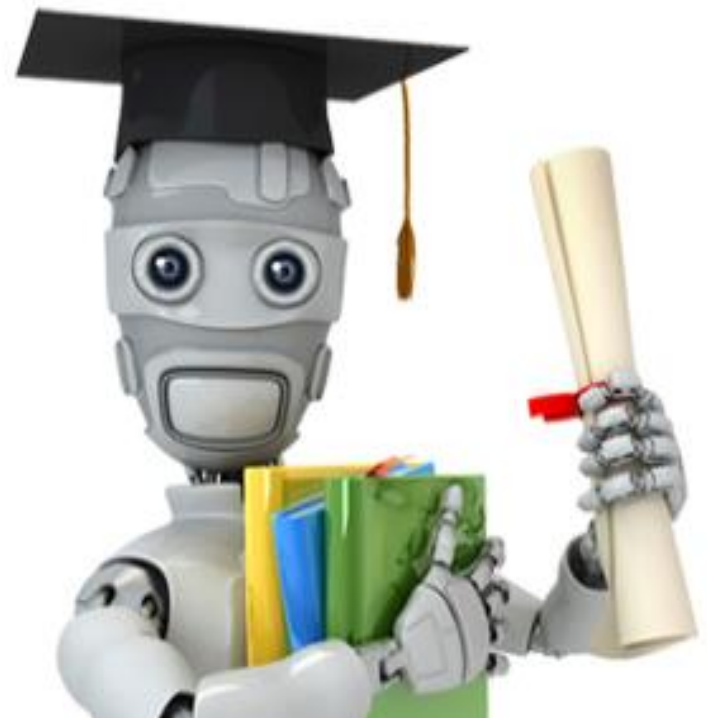
LinkedIn : [www.linkedin.com/in/andr -ch vez-a90078b9](https://www.linkedin.com/in/andr -ch vez-a90078b9).

# Agenda

- ✓ Introducción al Machine Learning.
- ✓ Supervised Learning.
- ✓ Unsupervised Learning.
- ✓ Reinforcement Learning.
- ✓ Deep Learning.
- ✓ Desarrollo de Algoritmos de Machine Learning.
- ✓ Implementación de Algoritmos de Machine Learning.
- ✓ Recursos.



# Introducción a los Algoritmos de Machine Learning



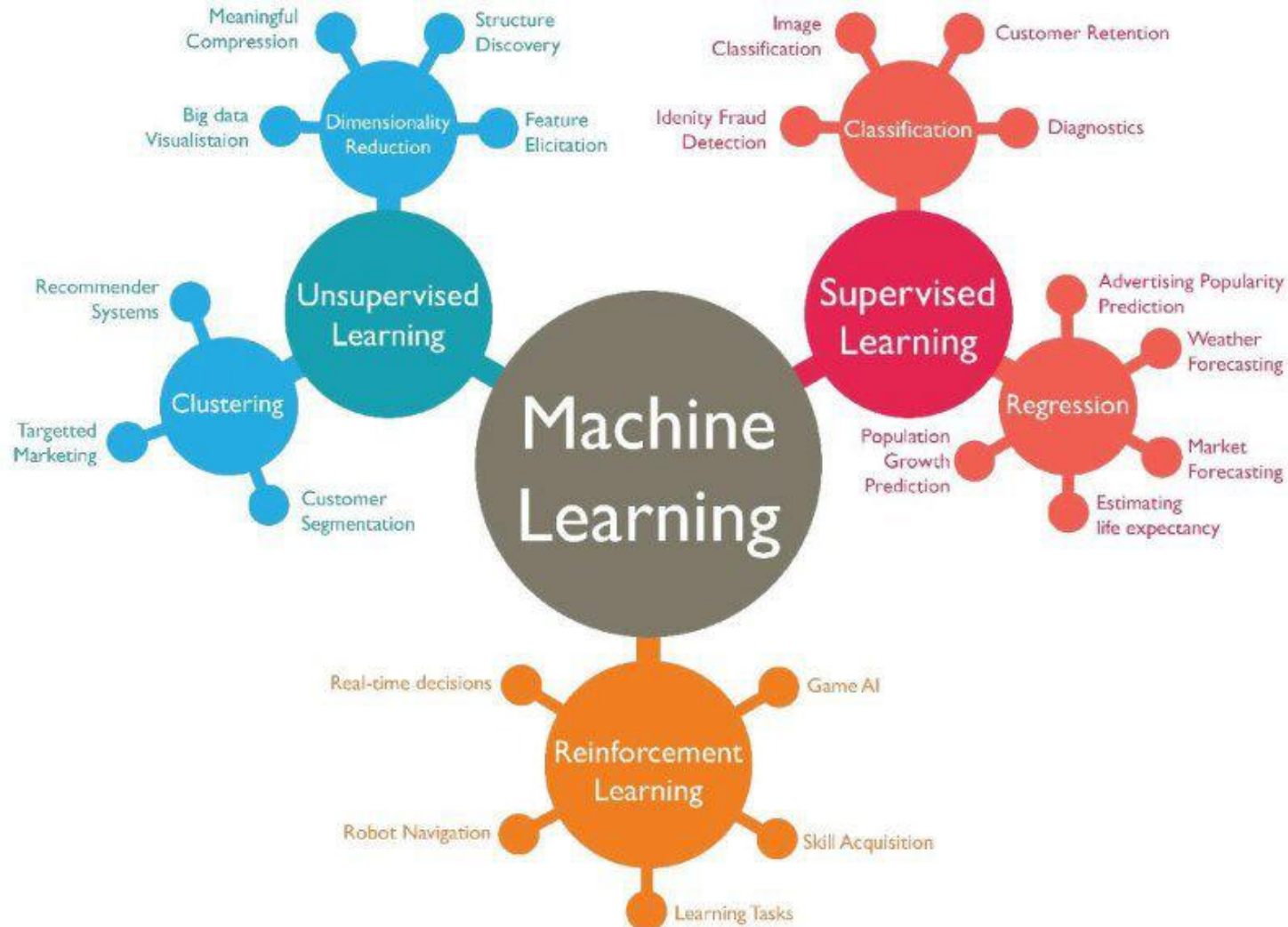
# ¿Qué es el Machine Learning (ML)?

- Rama de la inteligencia artificial que pretende que una máquina sea capaz de mejorar su actuación a la hora de resolver un problema mediante la adquisición de experiencia en una tarea determinada.
- **Multitud de algoritmos** con finalidades específicas.
- Ramas de Machine Learning:
  - ✓ Supervised Learning
  - ✓ Unsupervised Learning
  - ✓ Reinforcement Learning
  - ✓ Deep Learning

# Definiciones básicas

- **Conjunto de Datos (Data Set):** El total del conjunto de datos sobre los que queremos desarrollar un algoritmo de Machine Learning con el fin de obtener un modelo que lo represente lo mejor posible. Contendrá variables independientes y dependientes.
- **Variables Independientes (Features), (VI):** Aquellas columnas del Data Set que serán usadas por el algoritmo para generar un modelo que prediga lo mejor posible las variables dependientes.
- **Variables dependientes (Labels,Target), (VD):** Columna del data set que responde a una correlación de VI y que debe ser predicha por el futuro modelo
- **Conjunto de Datos de Entrenamiento (Training Set):** Subconjunto del Data Set que será utilizado para entrenar el modelo que se pretende generar.
- **Conjunto de Datos de Test (Test Set):** Subconjunto del data set que se le pasará al modelo una vez haya sido entrenado para comprobar, mediante el uso de diferentes métricas, sus indicadores más importantes de calidad.

# Introducción



# Supervised Learning (Modelos Supervisados)

- Se tiene una **variable objetivo** (Variable de Salida).
- Variables que ayudan a predecir a la variable de salida (Variables de entrada).
- Existe una dependencia de las variables de entrada con las variables de salida.





# Supervised Learning

- Género.



- Rangos de Edad.



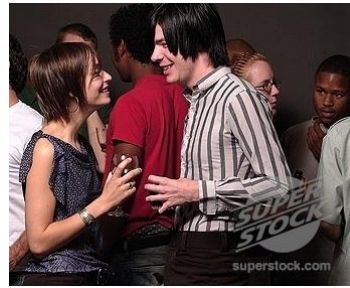
- Si Compra

- Ingresos.



- No Compra

- Estado Civil.





# Unsupervised Learning (Modelos No Supervisados)

- No hay una variable objetivo (Variable de Salida).
- No hay variables que ayudan a predecir a la variable de salida.



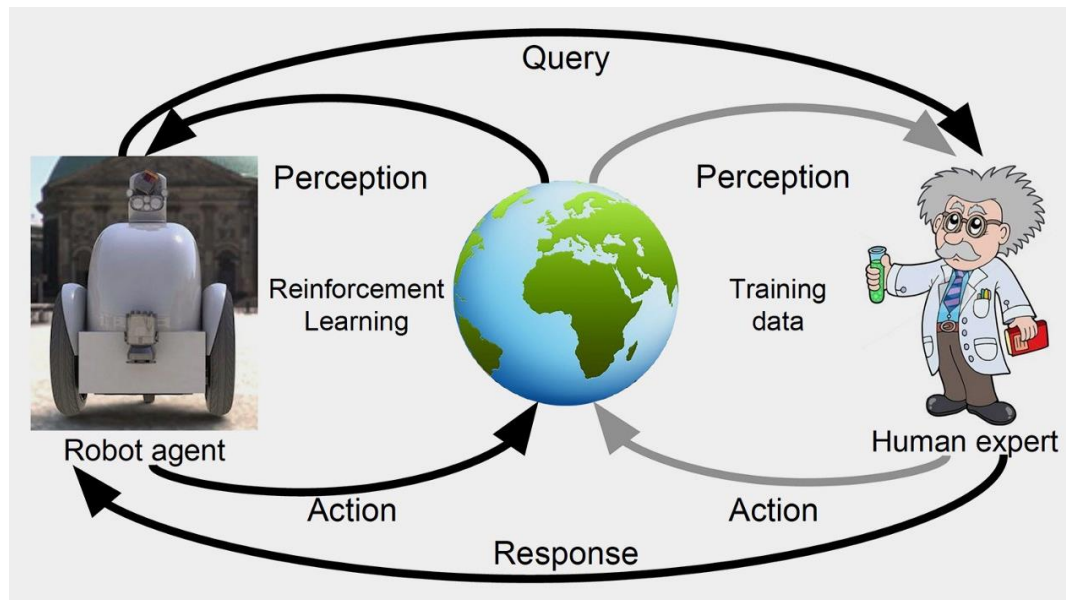
- Todas las variables tienen la misma importancia.
- Se busca la interdependencia de las variables.

# Modelos no Supervisados



# Reinforcement Learning (Aprendizaje por refuerzo)

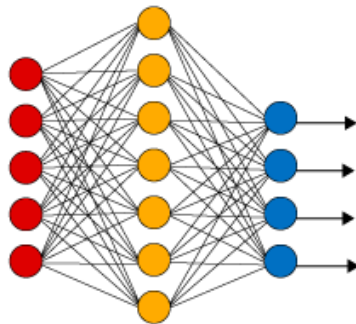
- El algoritmo de aprendizaje recibe un tipo de valoración acerca de la idoneidad de la respuesta dada.
- Cuando la decisión es correcta es muy parecido al aprendizaje supervisado, sin embargo difiere mucho cuando la decisión es incorrecta.



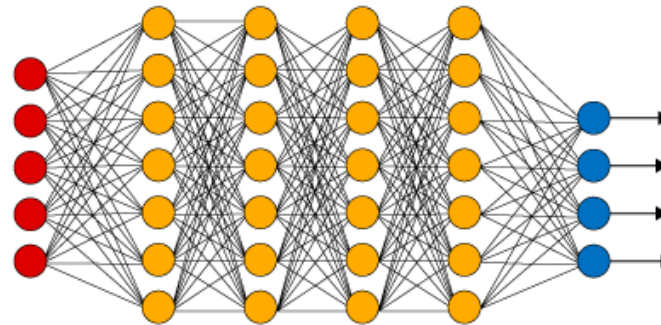
# Deep Learning (Aprendizaje Profundo)

- Es un conjunto de algoritmos de Machine Learning que intenta modelar abstracciones de alto nivel usando arquitecturas compuestas como redes neuronales profundas, redes neuronales convolucionales y redes de creencia profunda para resolver problemas como visión del computador, reconocimiento automático del habla, reconocimiento del audio y música, etc.

Simple Neural Network



Deep Learning Neural Network



● Input Layer

● Hidden Layer

● Output Layer

# Deep Learning (Aprendizaje Profundo)

## Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolution or Pool

Perceptron (P)



Feed Forward (FF)



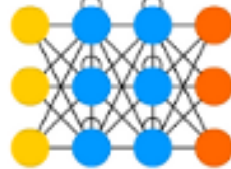
Radial Basis Network (RBF)



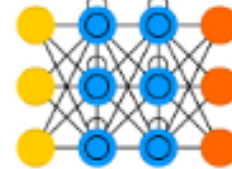
Deep Feed Forward (DFF)



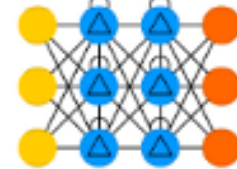
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



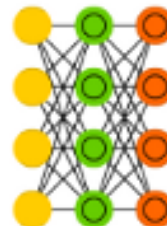
Gated Recurrent Unit (GRU)



Auto Encoder (AE)



Variational AE (VAE)



Denoising AE (DAE)



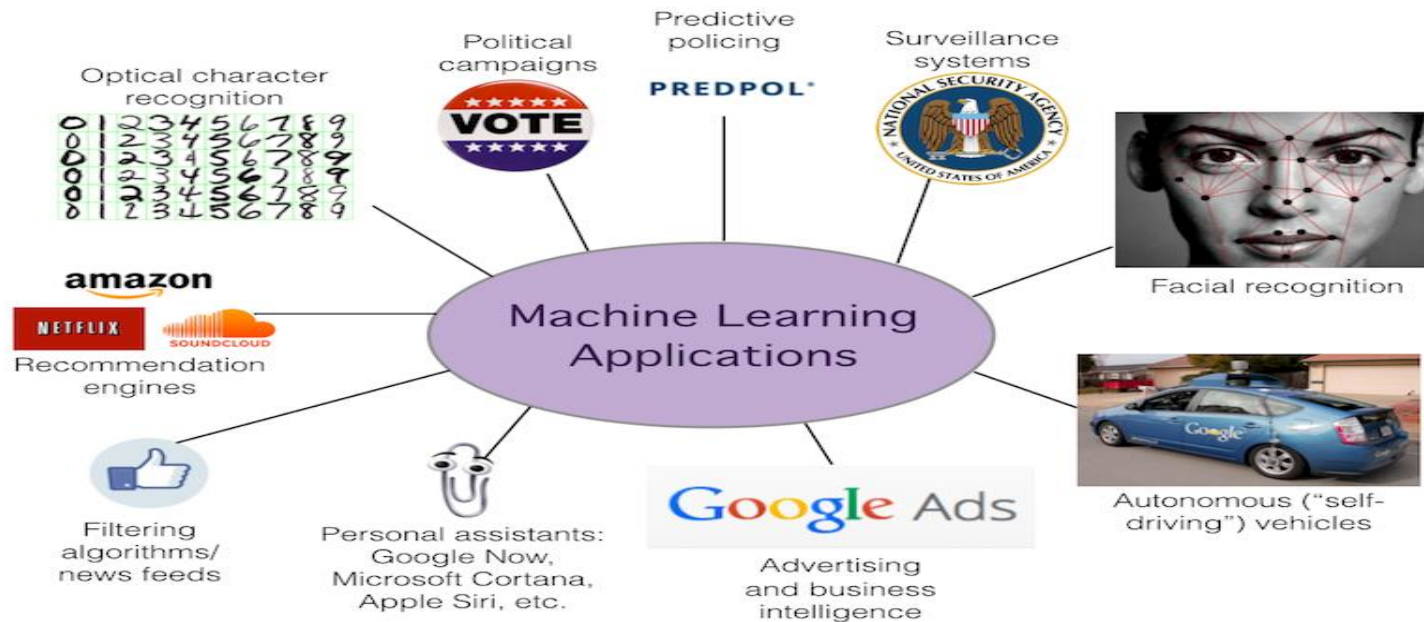
Sparse AE (SAE)



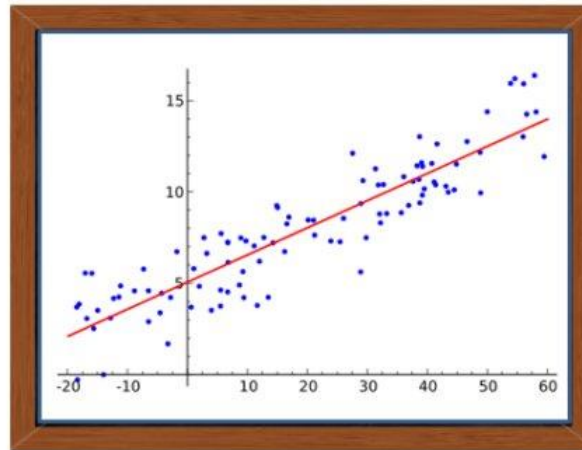


# Aplicaciones Machine Learning

- Sistemas de Recomendación.
- Detección de Spam.
- Natural Language Processing (NLP).
- Photo OCR (Optical Character Recognition).
- Visión Artificial.
- Diagnósticos médicos.
- Conducción Autónoma.
- AMD Ryzen ANN.



# Algoritmos Machine Learning: Regresión





# ALGORITMOS MACHINE LEARNING: REGRESIÓN

## Introducción

- Determinar la ecuación de regresión sirve para:
  - Describir de manera concisa la relación entre variables.
  - **Predecir los valores de una variable en función de la otra.**
- Veremos EXCLUSIVAMENTE relaciones lineales.
- La regresión lineal simple estudia la relación entre sólo dos variables (el caso de relación más sencillo posible).

# ALGORITMOS MACHINE LEARNING: REGRESIÓN

## Regresión lineal Simple y múltiple

- Para tratar este tipo de problemas se requiere expandir el análisis de regresión:

Regresión Lineal Simple



Regresión Lineal Múltiple

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$



$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

# ALGORITMOS MACHINE LEARNING: REGRESIÓN

## Interpretación del modelo de regresión lineal

$$\hat{Y} = 600 + 300X$$

- Supongamos que tenemos la ecuación de regresión, donde  $X$  es el número de años de experiencia profesional, e  $Y$  es el sueldo mensual.
  - ✓ Interpreta  $a$  y  $b$ .
  - ✓ Una persona con 3 años de experiencia laboral, ¿qué sueldo mensual tendrá? Interpreta el resultado.
  - ✓ Si una persona con 3 años de experiencia laboral tiene un sueldo mensual de 1700 €, ¿cuál será su error asociado?

# ALGORITMOS MACHINE LEARNING: REGRESIÓN

## Interpretación del modelo de regresión lineal

$$\hat{Y} = 600 + 300X$$

- ✓  $b=300 \rightarrow$  Cambio en  $Y$  por cada unidad de cambio en  $X$ . Por cada año de experiencia laboral, el sueldo mensual aumenta 300 €.
- ✓  $a=600 \rightarrow$  Valor medio de  $Y$  cuando  $X=0$ . Sueldo medio de aquellas personas sin experiencia laboral.

- ✓ Una persona con 3 años de experiencia laboral, ¿qué sueldo mensual tendrá? Interpreta el resultado.

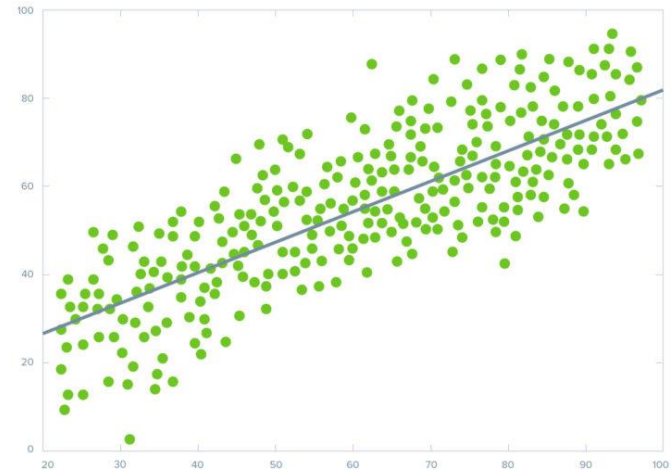
$$X = 3 \Rightarrow Y = 600 + 300 * 3 = 1500$$

$$\hat{Y} = 1500$$

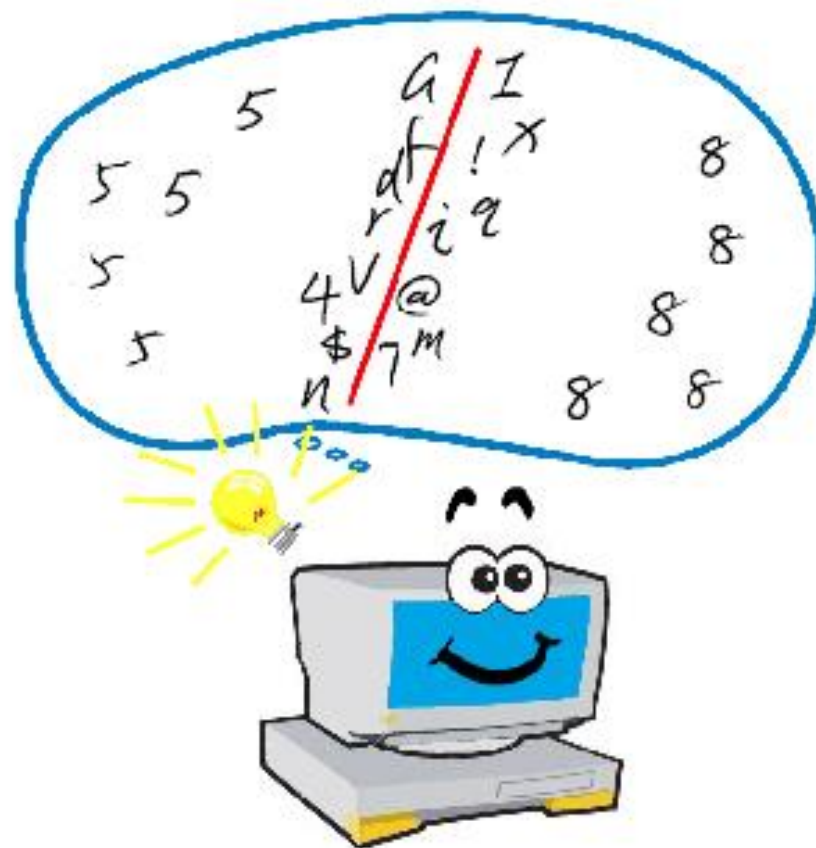
- ✓  $\rightarrow$  Valor promedio previsto para todos los sujetos que han obtenido en la variable  $X$  un valor de  $X_i$ . Las personas con 3 años de experiencia tienen un sueldo promedio de 1500 €.

# Algoritmos de regresión

- ✓ **Simple Linear Regression.**
- ✓ **Multiple Linear Regression.**
- ✓ **Polynomial Regression.**
- ✓ **Support Vector Regression (SVR).**
- ✓ **Decision Tree Regression.**
- ✓ **Random Forest Regression.**
- ✓ **XGBoost, LightGBM Regression.**

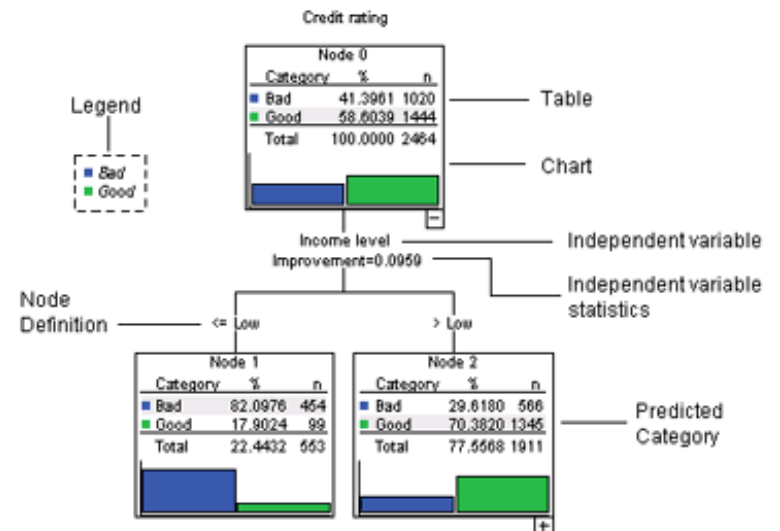


# Algoritmos Machine Learning: Clasificaciòn



# ALGORITMOS MACHINE LEARNING: CLASIFICACIÓN

## Árboles de Clasificación

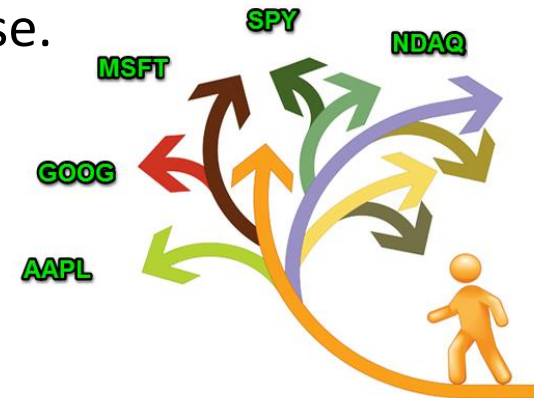




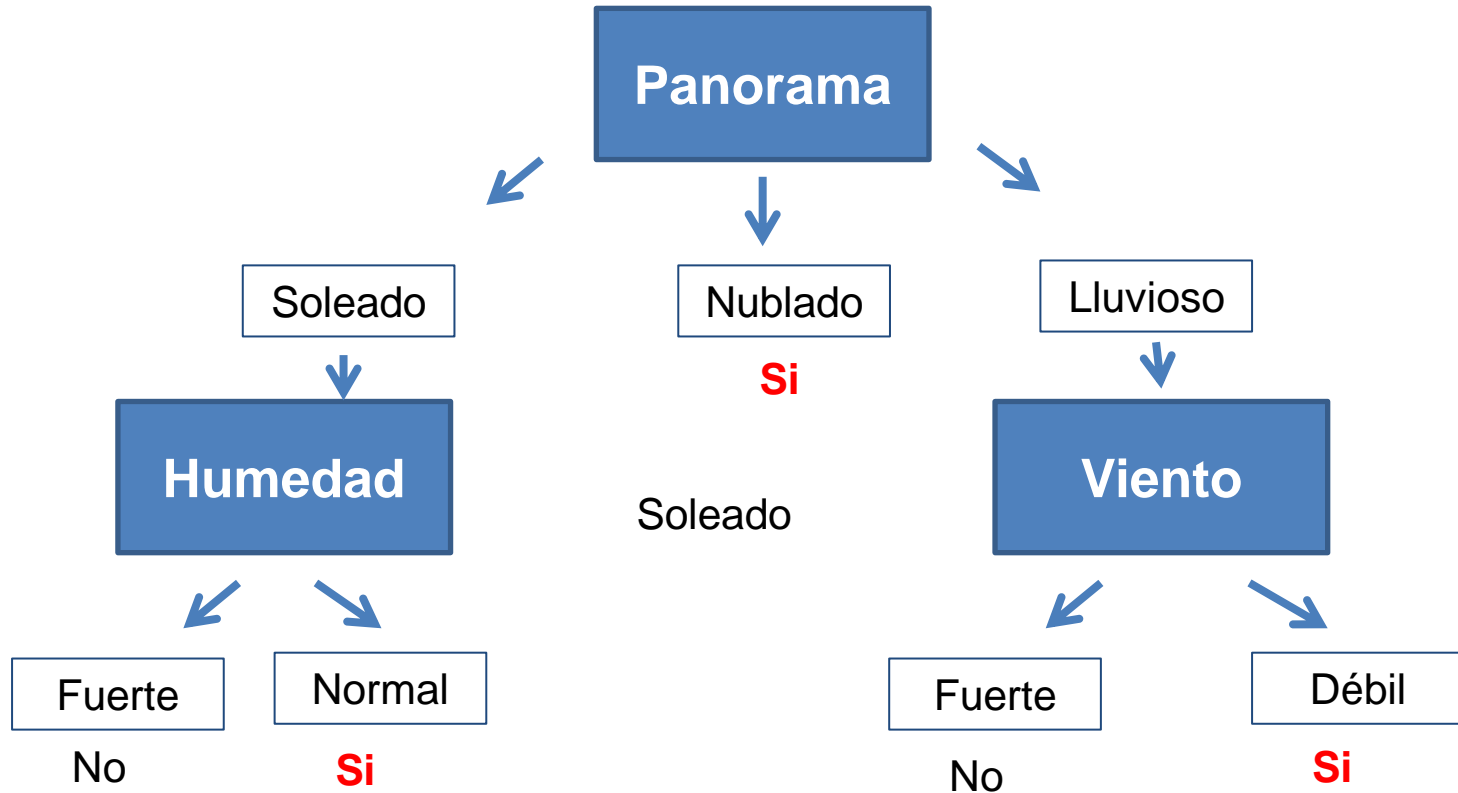
# ALGORITMOS MACHINE LEARNING: CLASIFICACIÓN

## Árboles de Clasificación

- Entrada:
- Objetos caracterizables mediante propiedades.
- Salida:
  - En árboles de decisión: una decisión (sí o no).
  - En árboles de clasificación: una clase.
- Conjunto de reglas.

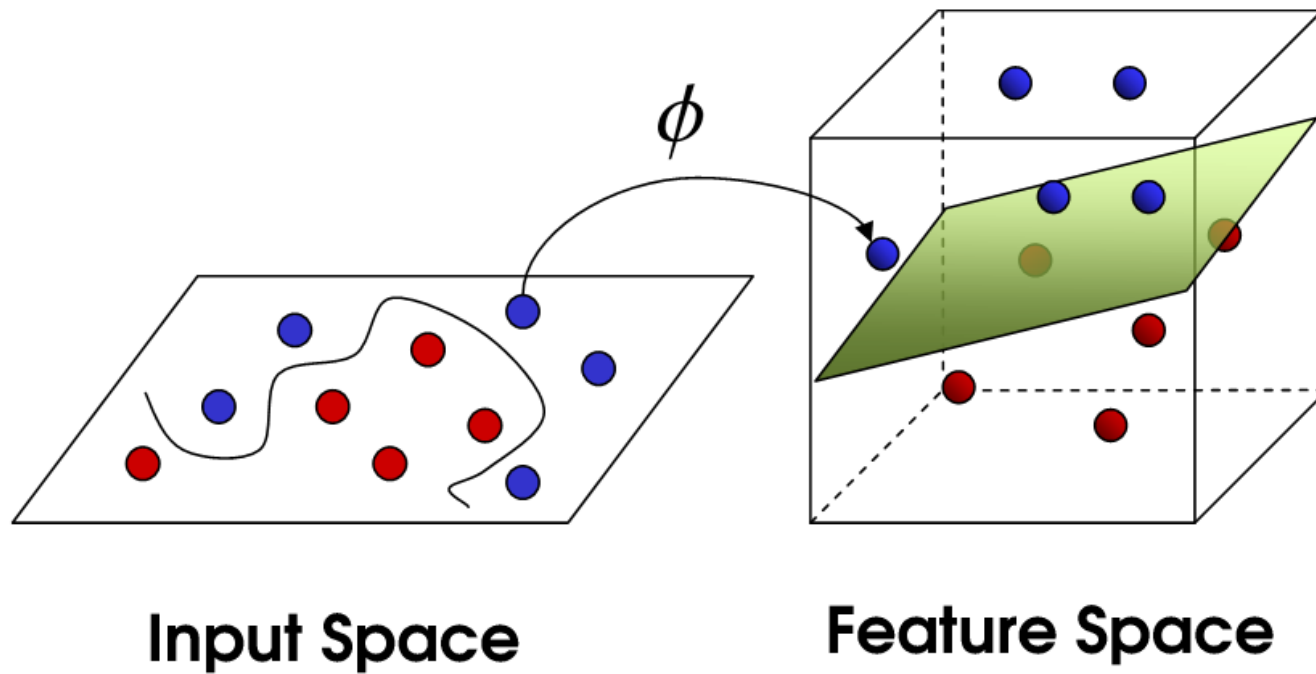


# Juego Tennis ?



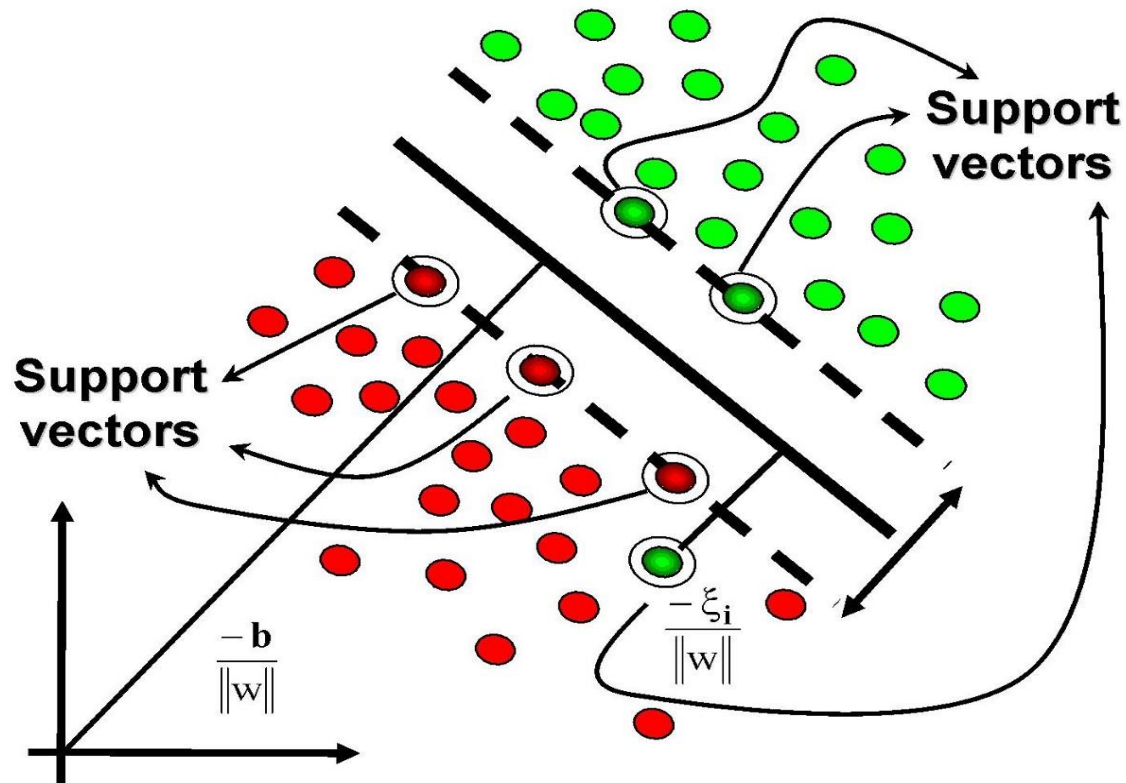
# ALGORITMOS MACHINE LEARNING: CLASIFICACIÓN

## SVM



# ALGORITMOS MACHINE LEARNING: CLASIFICACIÓN

## SVM



# Algoritmos de clasificación

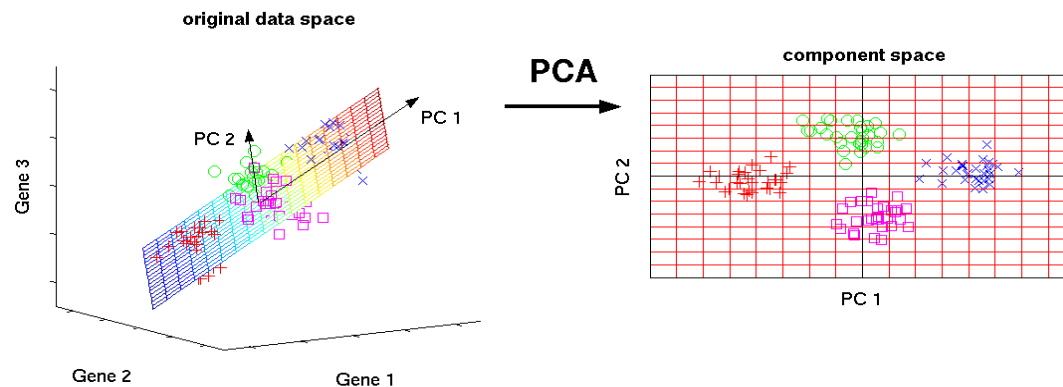
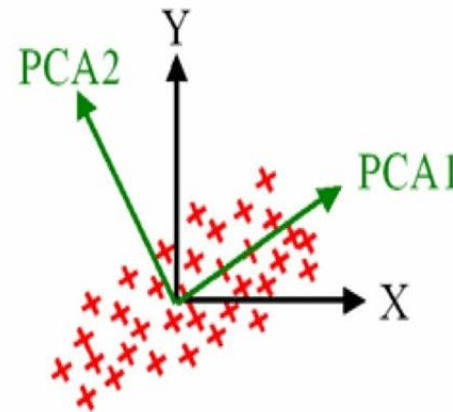
- ✓ Logistic Regression
- ✓ **Support Vector Machine (SVM)**
- ✓ Naïve Bayes
- ✓ **Decision Tree Classification**
- ✓ Random Forest Classification
- ✓ LGBM , CatBoost y XGBoost.



# Algoritmos machine learning: Reducción de Dimensión

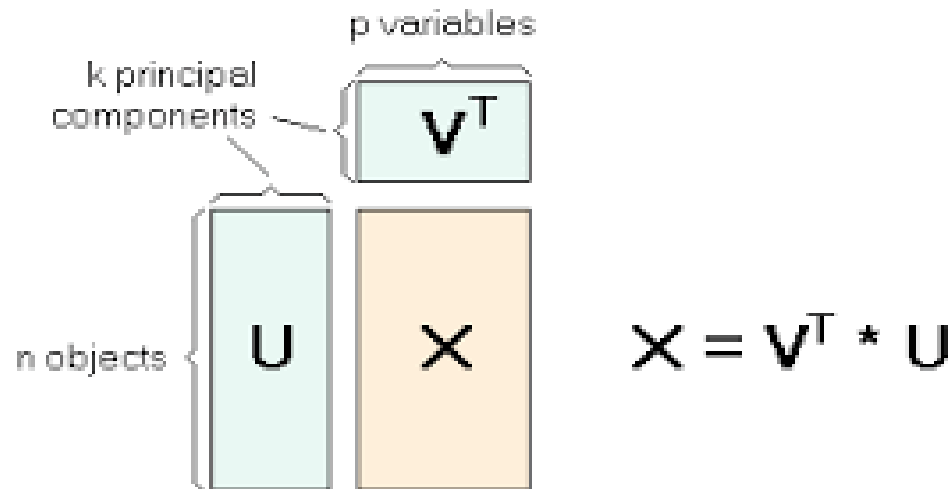


- **Karl Pearson**



## ALGORITMOS MACHINE LEARNING: REDUCCIÓN DE DIMENSIÓN

- **Objetivo:** Dada una matriz de datos de dimensiones  $n \times p$  que representa los valores de  $p$  variables en  $n$  individuos, investigar si es posible representar los individuos mediante  $k$  variables ( $k < p$ ) con poca (o ninguna si es posible) pérdida de información.



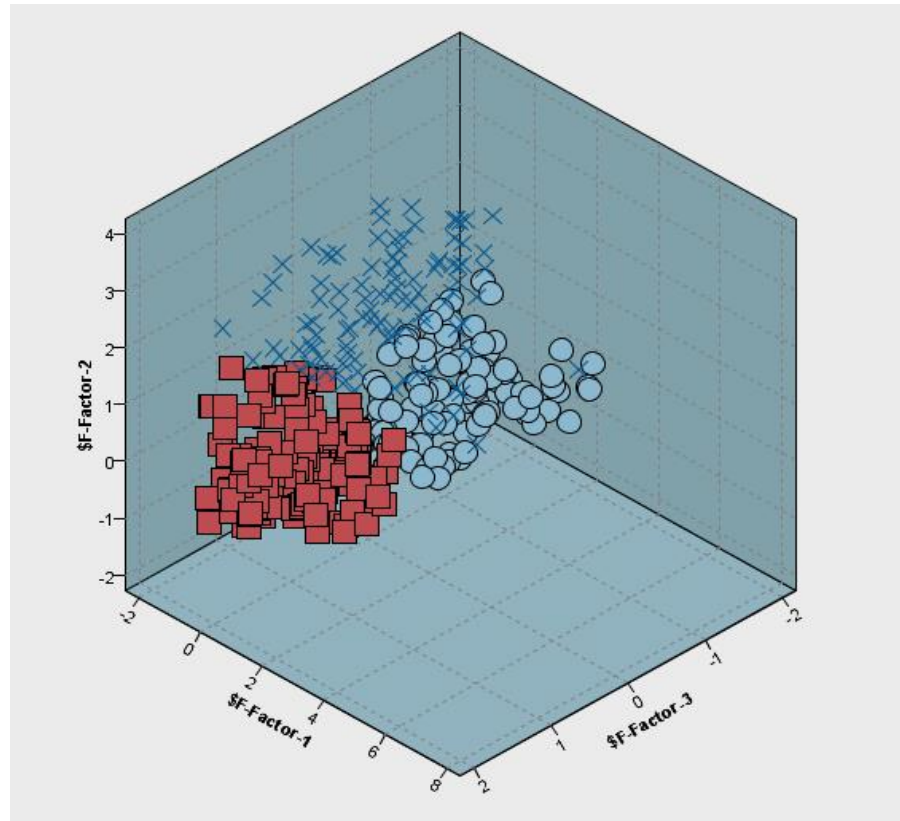


# Algoritmos machine learning: Clustering



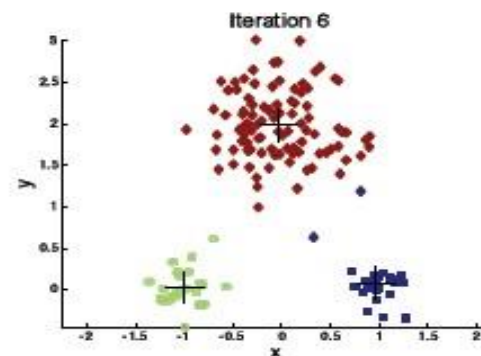
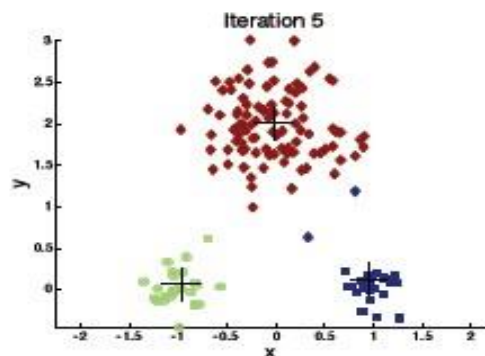
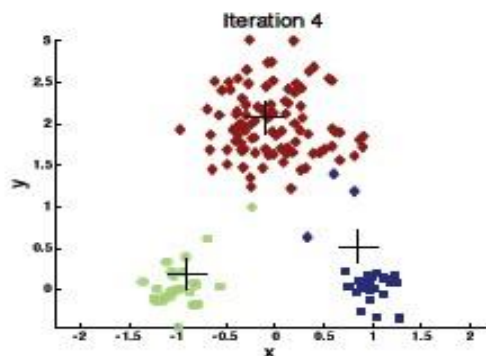
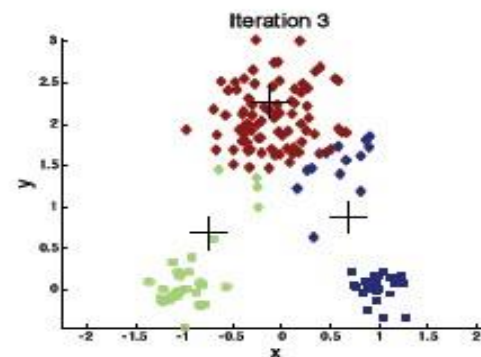
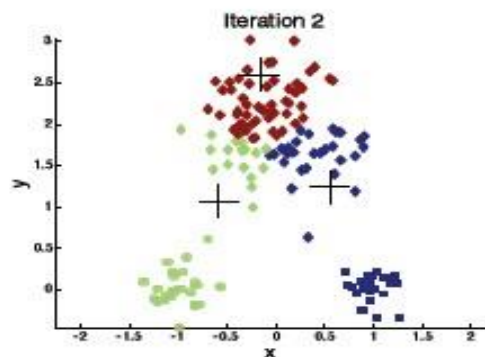
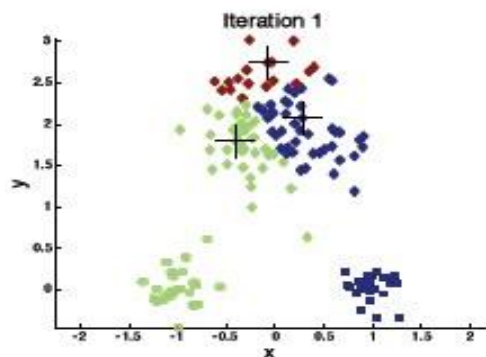
## ALGORITMOS MACHINE LEARNING: CLUSTERING

# ALGORITMO DE K - MEANS



## ALGORITMOS MACHINE LEARNING: CLUSTERING

# Método de K - Means



## ALGORITMOS MACHINE LEARNING: CLUSTERING

### Teorema: Igualdad de Fisher

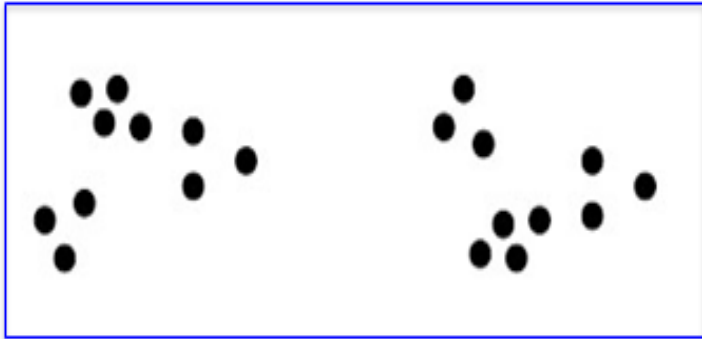
*Inercia total = Inercia Inter - clases*

*+*

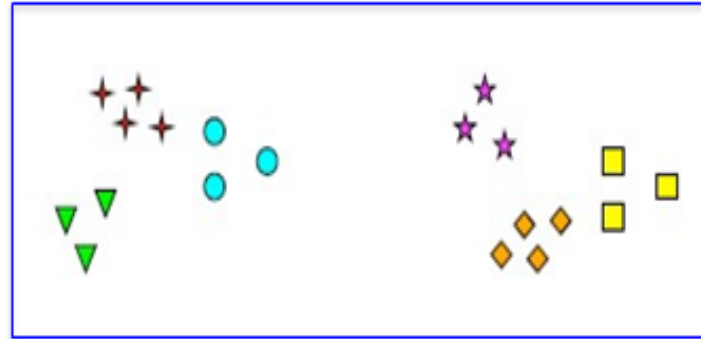
*Inercia Intra-clases*

$$I = B(P) + W(P)$$

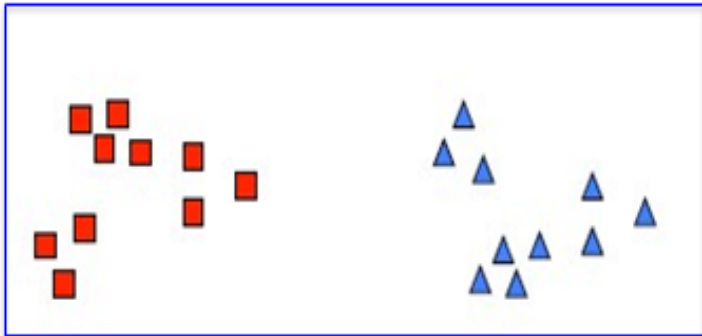
# ¿ Cuántos clústeres?



Datos originales



6 clústeres

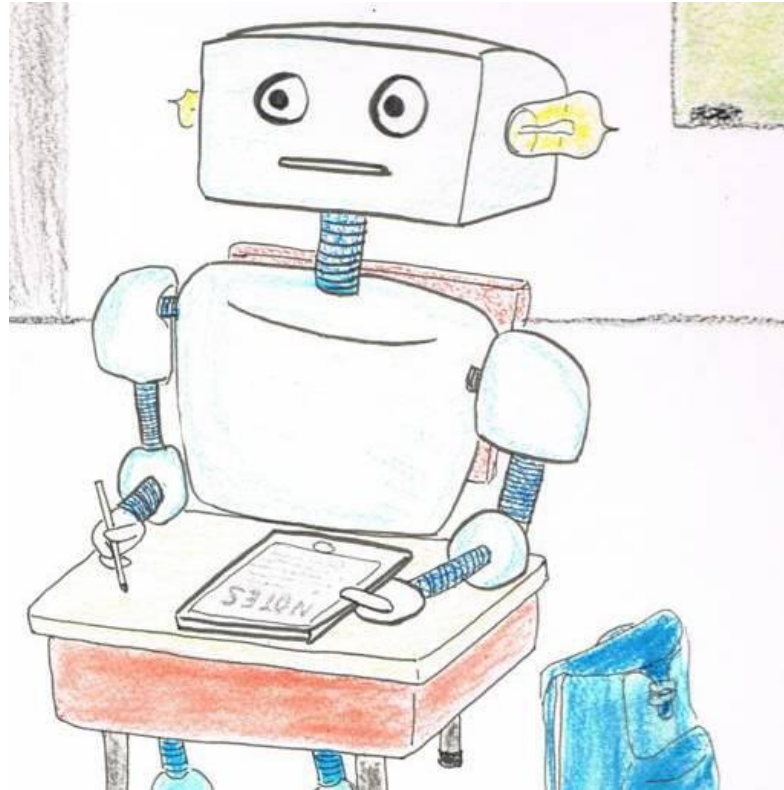


2 clústeres



4 clústeres

# Evaluando un Algoritmo de Machine Learning



# Evaluando un Algoritmo de Machine Learning

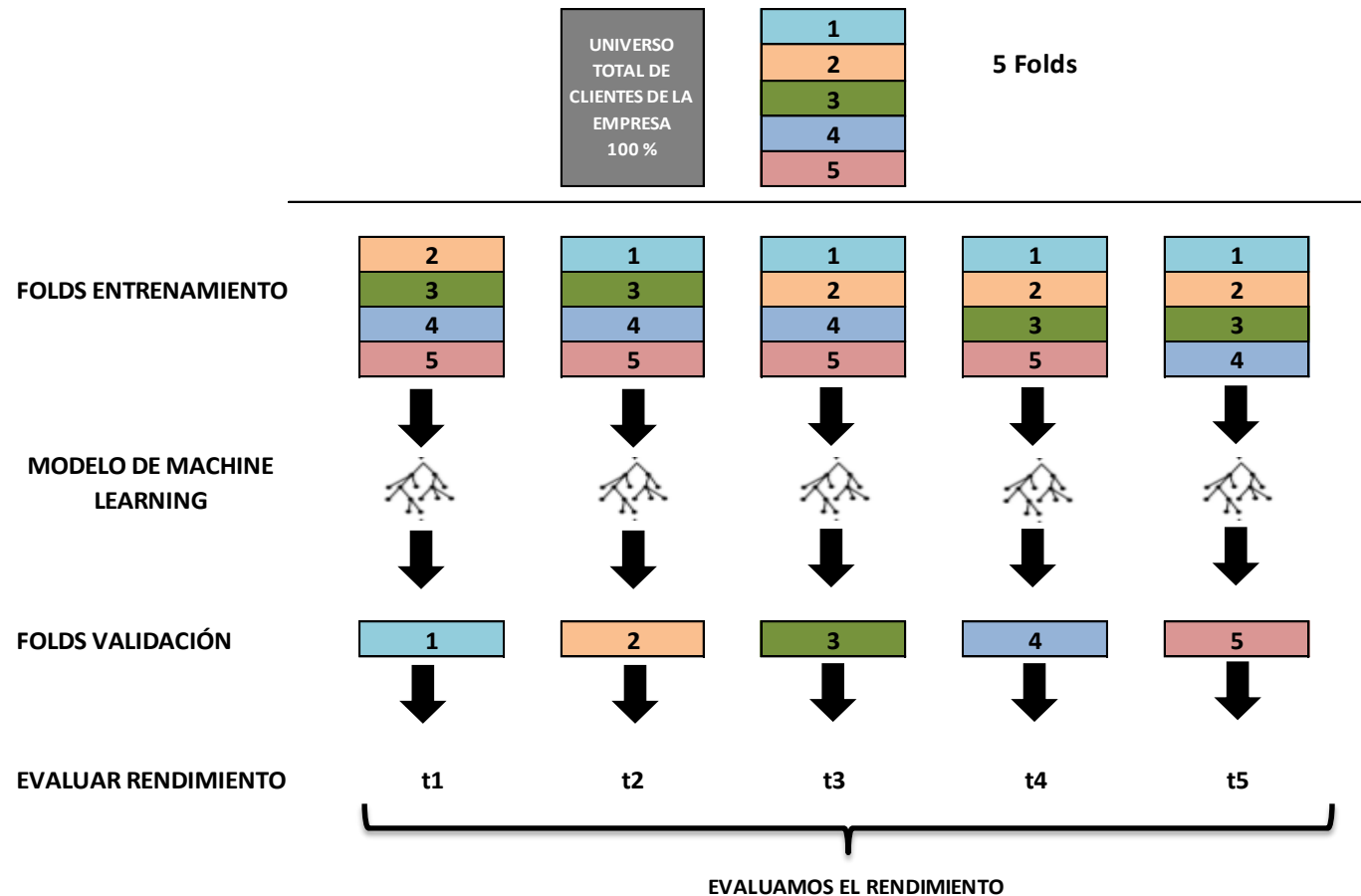
## MUESTRA DE ENTRENAMIENTO Y VALIDACIÓN





# Evaluando un Algoritmo de Machine Learning

## VALIDACIÓN CRUZADA



# Evaluando un Algoritmo de Machine Learning

## MATRIZ DE CONFUSIÓN Y MATRIZ DE COSTOS

MATRIZ DE CONFUSIÓN		PREDICCIÓN	
		NO MOROSOS	MOROSOS
REALIDAD	NO MOROSOS	DECISIÓN CORRECTA VN	<b>FP</b>
	MOROSOS	<b>FN</b>	DECISIÓN CORRECTA VP

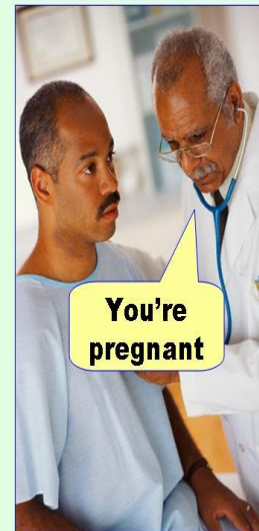
$$\text{PRECISIÓN} = (VN + VP) / (VN + VP + FP + FN)$$

$$\text{SENSIBILIDAD} = VP / (VP + FN)$$

$$\text{ESPECIFICIDAD} = VN / (VN + FP)$$

$$\text{F-SCORE} = 2 * ((VP / (VP + FP)) * (VP / (VP + FN))) / ((VP / (VP + FP)) + (VP / (VP + FN)))$$

**Type I error**  
(false positive)



**Type II error**  
(false negative)



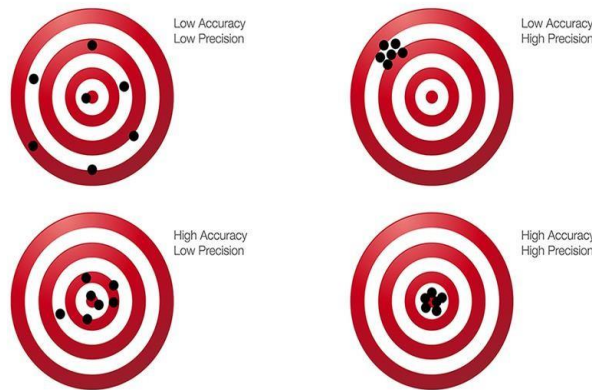
# Aplicación de Machine Learning

## Caso práctico: Clasificación del cáncer

- Casos en los que el nº de ejemplos negativos es mucho mayor que el de ejemplos positivos
- Ejemplo:

- Modelo regresión logística
  - $y = 1$       *cáncer*
  - $y = 0$       *no cáncer*
- Se tiene un 1 % de error en el set de test (99 % de diagnósticos correctos)
- Sólo el 0,5 % de los pacientes tiene cáncer

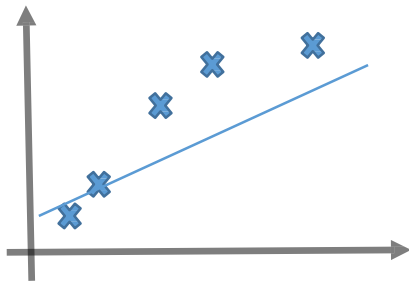
## Exactitud vs. Precisión (Accuracy vs. Precision)



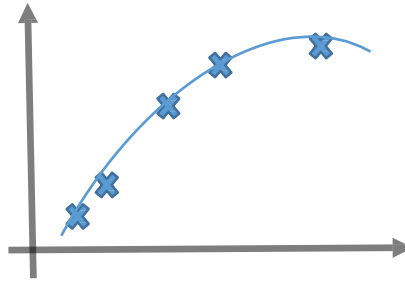
# El problema del sobreajuste y Sobregeneralización

- La sobregeneralización (Underfitting) se da cuando un modelo es demasiado simple y no se ajusta a los datos de entrenamiento.
- El sobreajuste (Overfitting) aparece cuando un modelo es muy complejo y se ajusta demasiado bien a los ejemplos de entrenamiento pero mal a los de test (no generaliza)

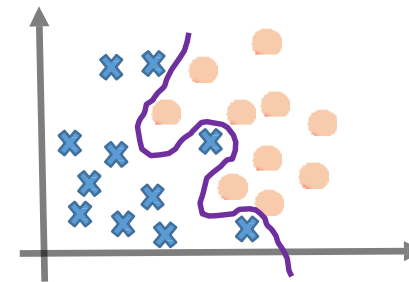
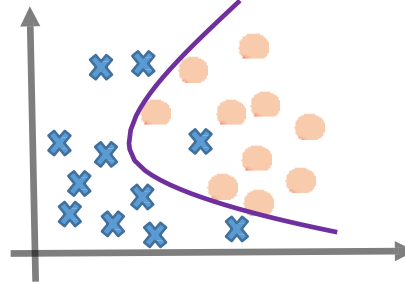
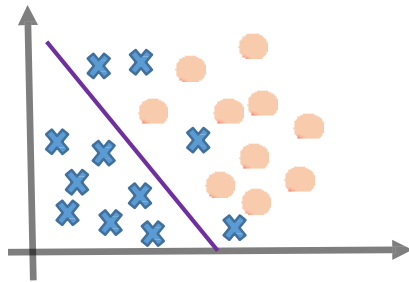
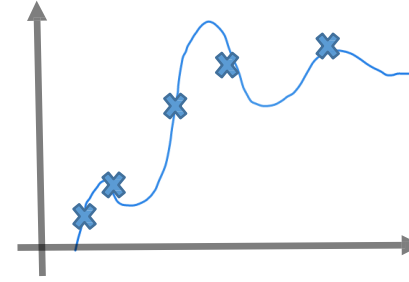
“Underfitting”, “High Bias”



OK



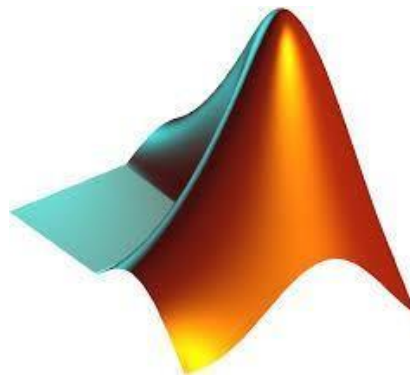
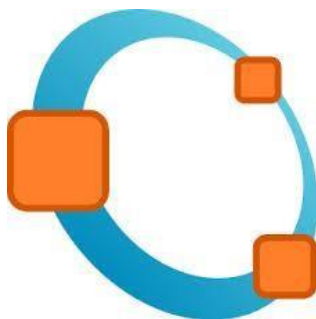
“Overfitting”, “High variance”



# Panorama tecnológico : Softwares Machine Learning

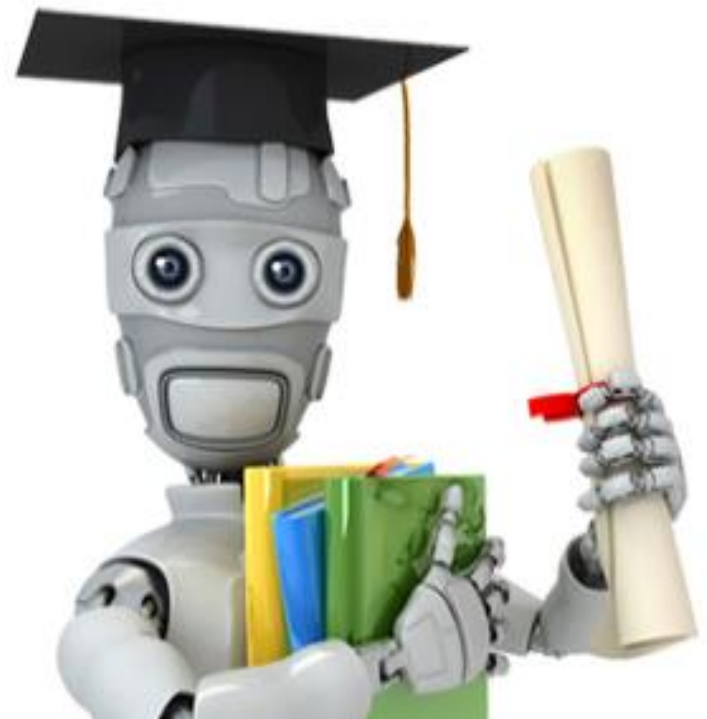


SPSS Modeler

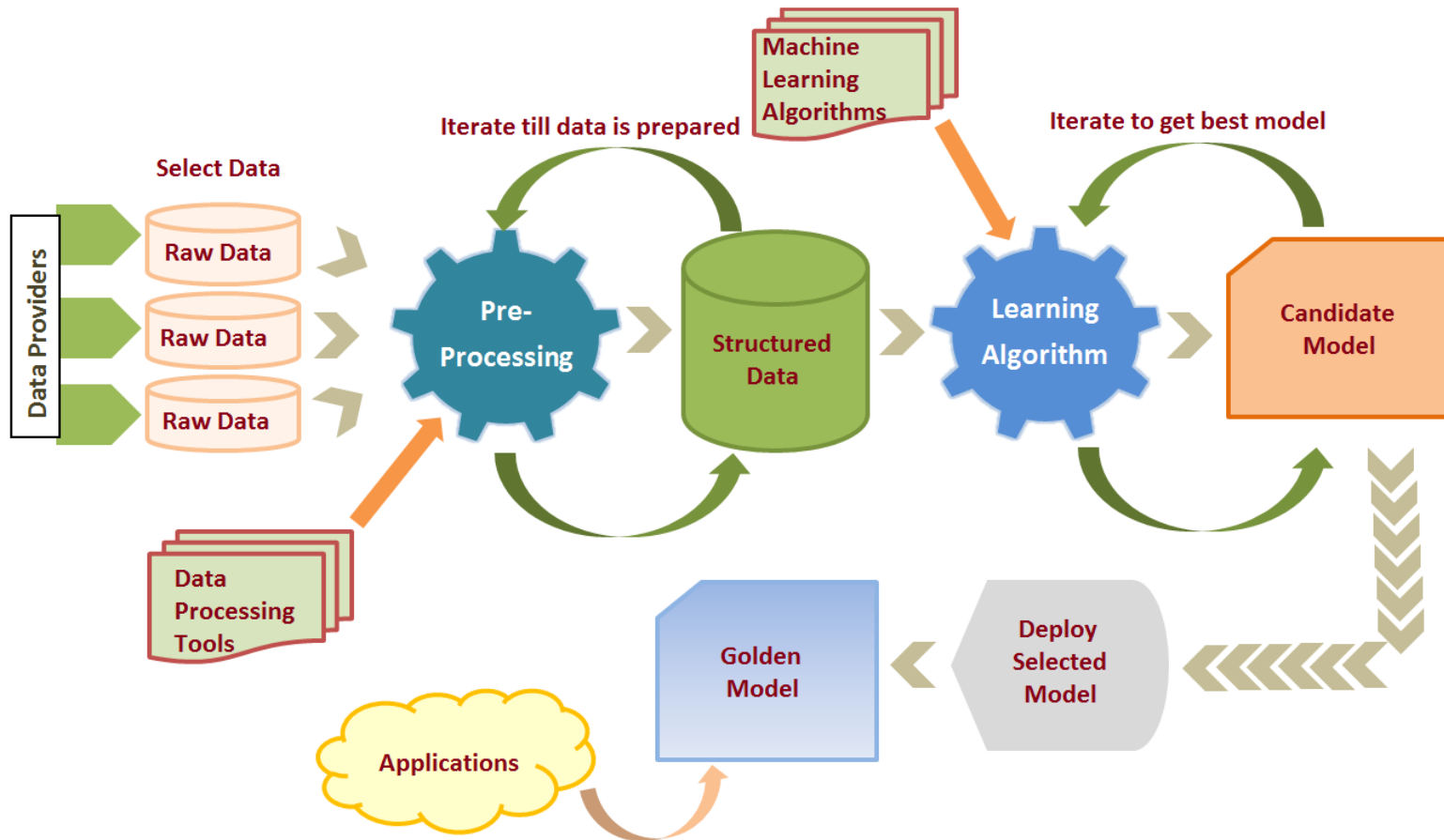


# **Desarrollo de Algoritmos de Machine Learning**

**70 %**



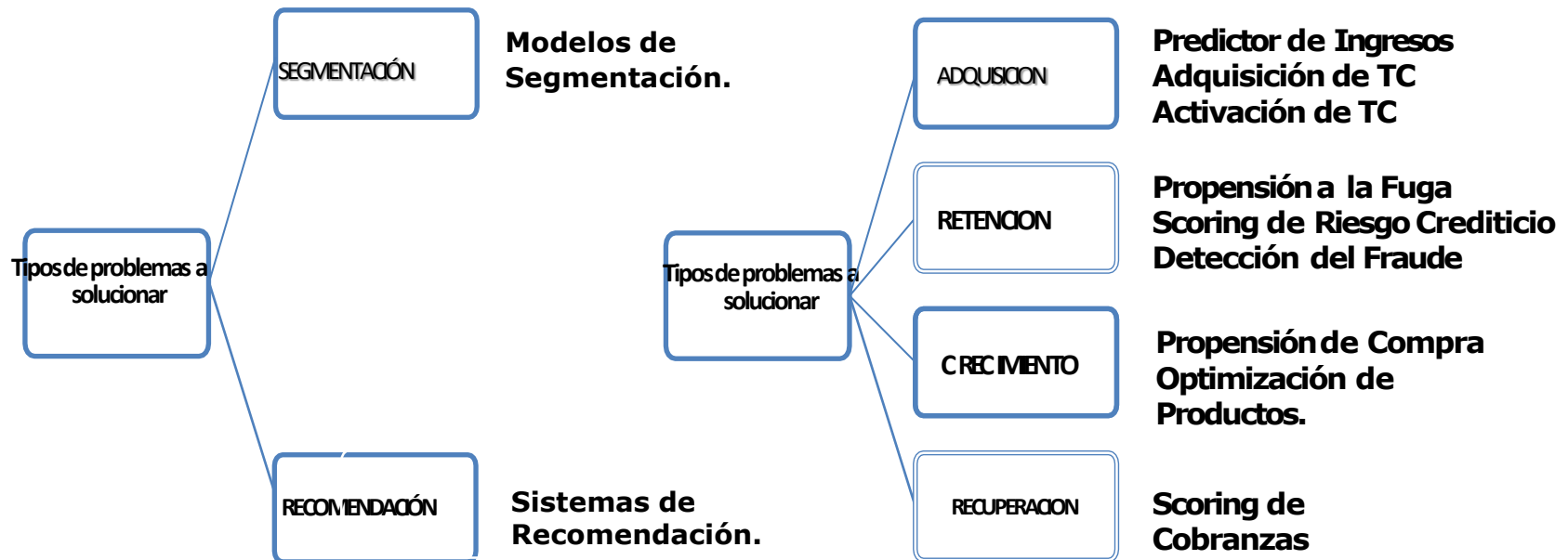
# Machine learning en los negocios



# ENTENDIMIENTO DEL PROBLEMA

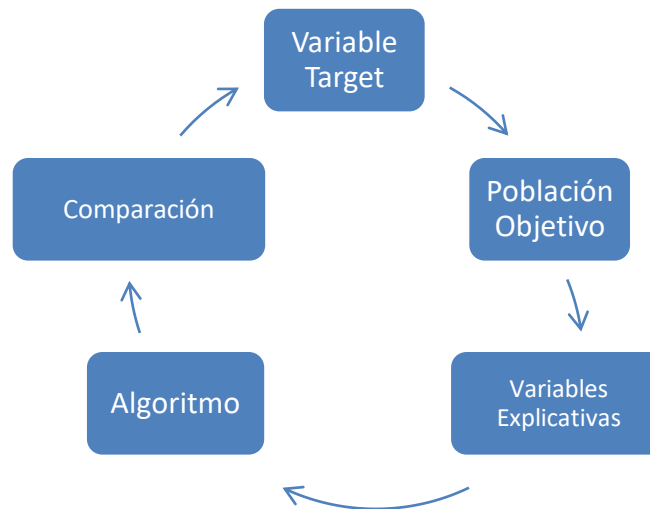
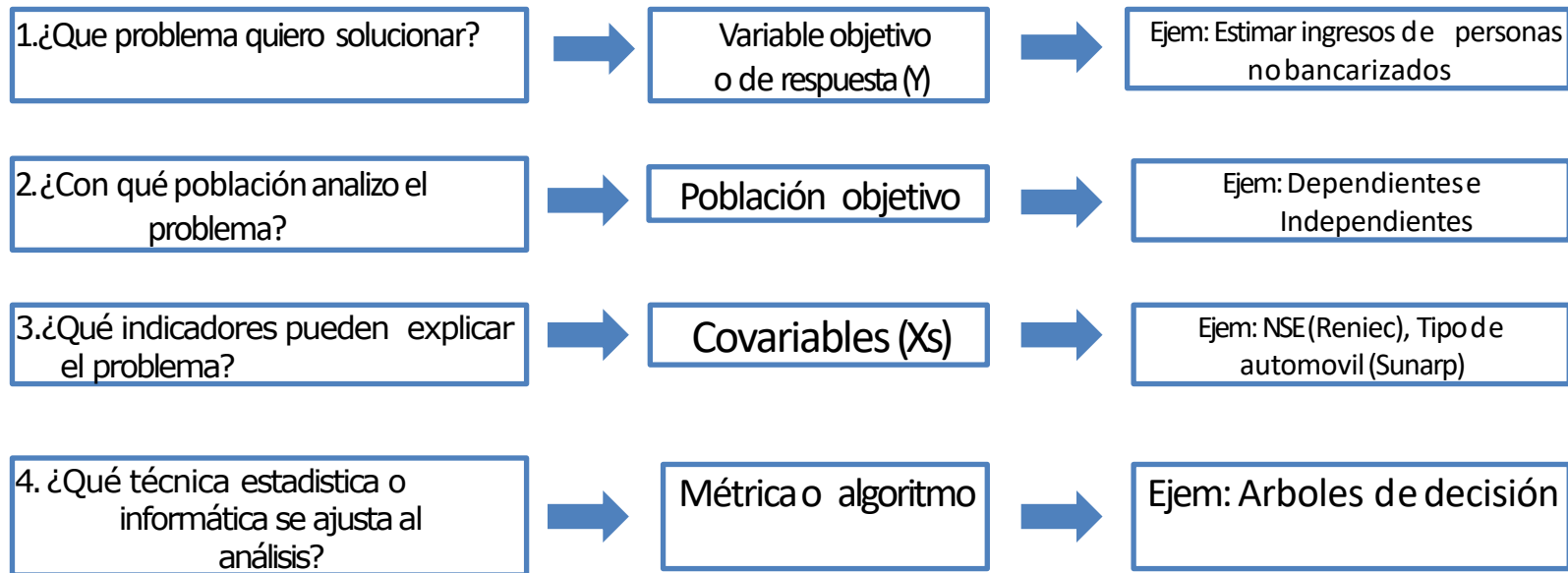
## PROPÓSITO DEL ANÁLISIS

Descubrir eventos o resultados futuros en base al conocimiento previo de los datos, utilizando para ello métodos estadísticos, matemáticos, computacionales y de base de datos, así como de la aplicación de los algoritmos de machine learning. En cualquier negocio el éxito depende de:

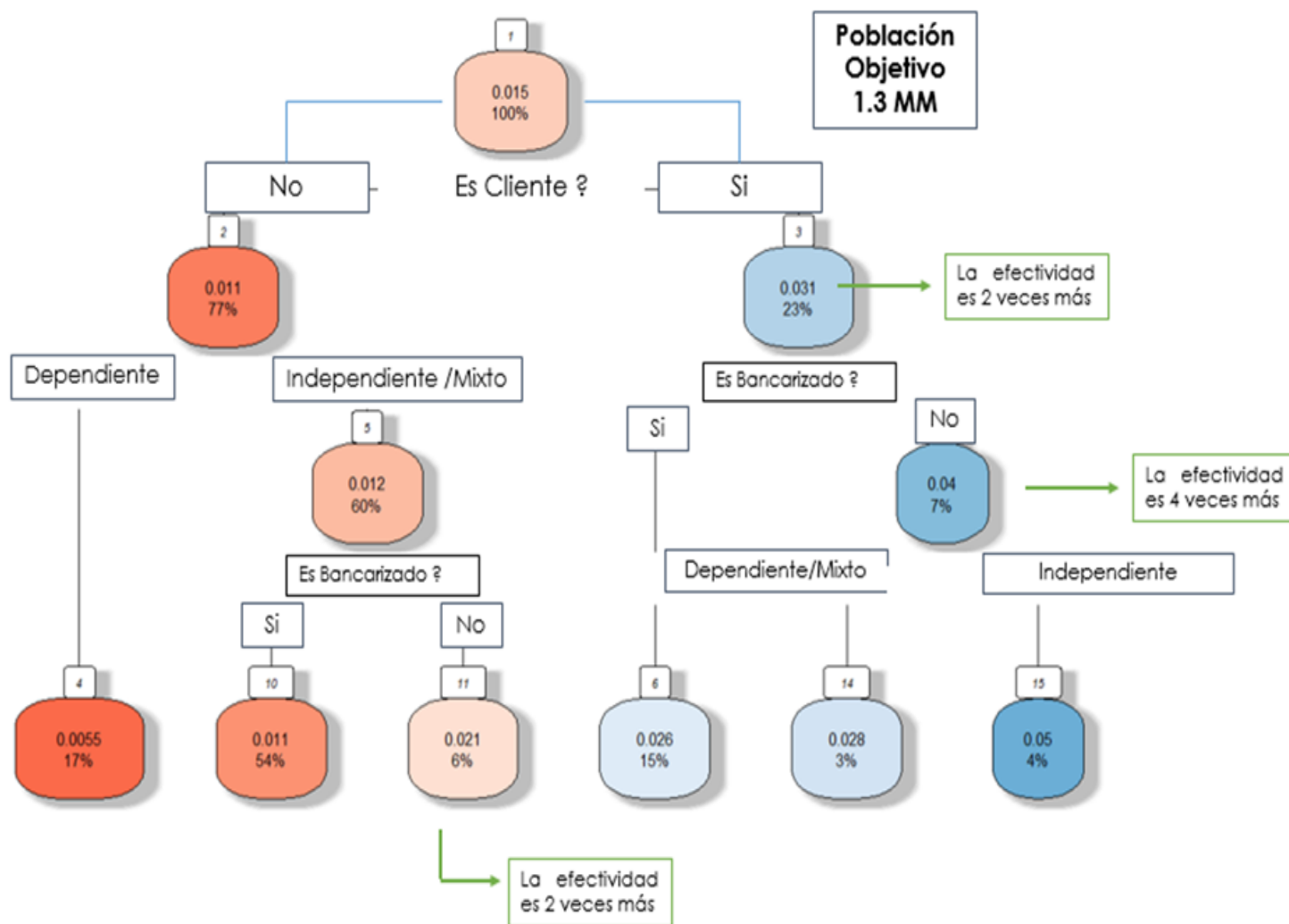




# ENTENDIMIENTO DEL PROBLEMA Y POSIBLE SOLUCIÓN



# SEGMENTACIÓN DE LA PO



# CRITERIOS DE INCLUSIÓN Y EXCLUSIÓN DE LA INFORMACIÓN

## 1. Data Original

```
100100011101000000101000110111010110
100100111101110000001111100110100100
100001101101111101010011100001101001
111111010000110111001010111100001011
1100111110111111100100001110110110
010000110100110110000110000100010000
010101110011001111011001110100010111
001000010101100101000001000010011110
0111010011111100101110101010111100
1000100001011000101011010111000101
01001000010010101110011100001010000
0101100000100111010100101110110001
011011111010111100010100010100010000
011010011011011010001000101111001101
00010100000110011000110010001001110
10010101010001001110010101011111101
```

Criterios de  
exclusión

## 2. Data de Estudio

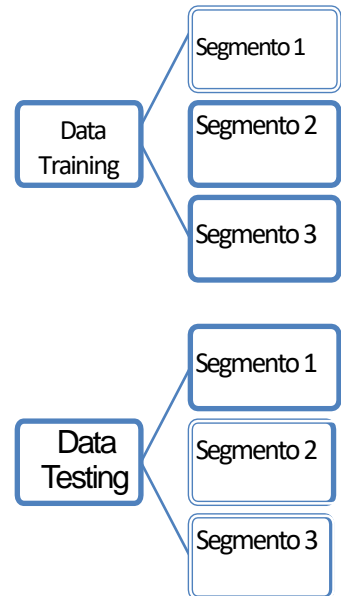
```
100100011101000000101000110111010110
100100111101110000001111100110100100
100001101101111101010011100001101001
111111010000110111001010111100001011
1100111110111111100100001110110110
010000110100110110000110000100010000
010101110011001111011001110100010111
001000010101100101000001000010011110
0111010011111100101110101010111100
1000100001011000101011010111000101
01001000010010101110011100001010000
01001000010010101110011100001010000
```

Training

Testing

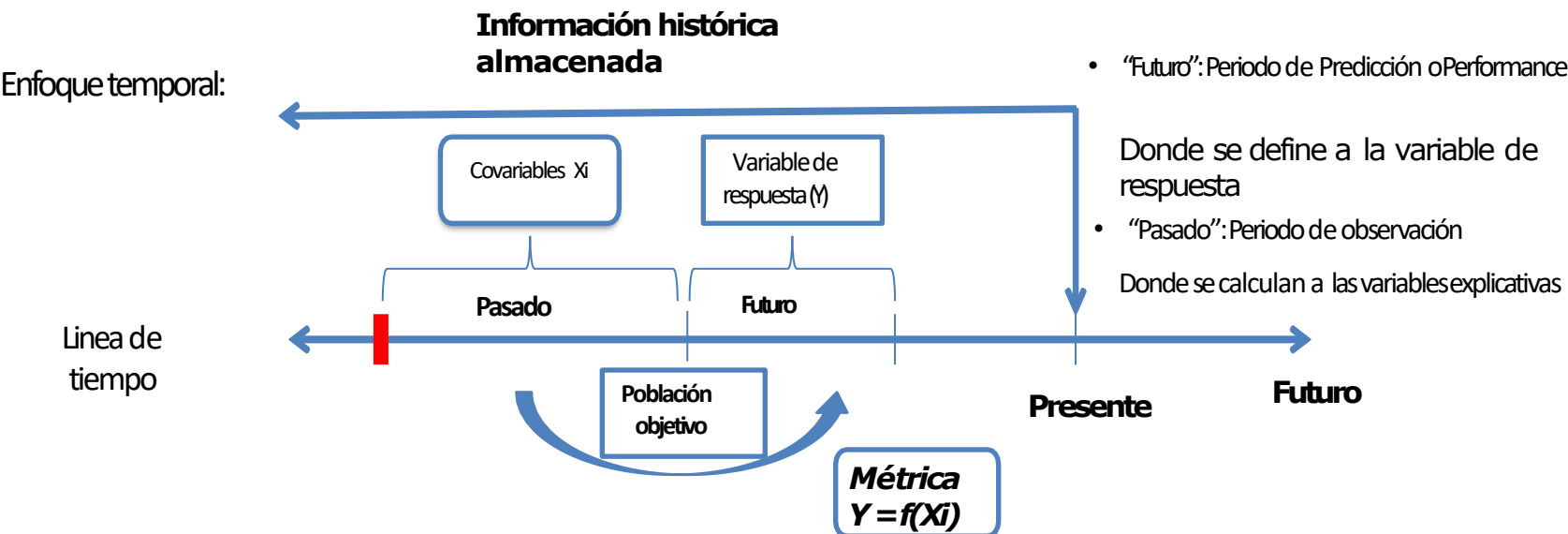
Información pulida y lista  
para el proceso de  
modelación

## 3. Segmentación



- ✓ Tienen poca información histórica
- ✓ Son riesgosos para el negocio
- ✓ Están fuera de las políticas o de las estrategias del negocio
- ✓ Son datos erróneos
- ✓ Son casos especiales que no se volverán a recopilar
- ✓ Ciclos económicos
- ✓ Temas regulatorios / legales

# DEFINICIÓN DE LA VARIABLE TARGET Y HORIZONTE TEMPORAL



Enfoque matricial:

ID	Segment_Target	Var_Target	Var_X1	Var_X2	Var_X3	Var_X4	Var_X5	Var_X6
1	Segment 1	1	-0.243257655	216	952.4800	1	4	3
2	Segment 2	1	1.696358794	191	633.4949	0	7	2
3	Segment 3	1	0.561226988	192	637.5107	0	6	3
4	Segment 1	1	-1.673888687	205	927.2513	0	8	3
5	Segment 2	0	-0.315746538	200	988.0877	0	2	3
6	Segment 3	0	0.402197729	201	927.5218	1	6	2
7	Segment 1	1	0.668736379	202	582.0028	0	6	2
8	Segment 2	1	1.489475004	197	701.1748	0	6	2
9	Segment 3	0	0.308647509	201	526.3747	0	8	4
10	Segment 1	1	0.090616380	189	989.2571	0	7	4
11	Segment 2	1	0.081223506	200	789.0298	0	8	2
12	Segment 3	1	-0.443663814	207	937.3809	0	2	2
13	Segment 1	1	-1.416088194	220	819.6118	0	9	1
14	Segment 2	1	-0.316298576	187	995.7736	1	2	5

Población objetivo

Variable de respuesta (Y)

Covariables Xi

Métrica  
 $Var\_Target = f(Var\_X1, Var\_X2, Var\_X3, Var\_X4, Var\_X5, Var\_X6)$

# Periodicidad , Cosechas y Matrices de Transición

		Ene	Feb	Mar	Abr	May	Jun	Jul	Ago	Sep	Oct	Nov	Dic
Test	1												
	2												
Train	1												
	2												
	3												



Información  
histórica



Ocurrencia de la  
Target

## DEFINICIÓN Y CREACIÓN DE DRIVERS

Las variables a seleccionar para la solución del problema propuesto deben tener **sentido para el negocio**. En otras palabras al seleccionarlasy se espera que estén correlacionadas con la variable de respuesta del modelo. La transformación tiene como propósito optimizar el aporte de las  $X_i$  en el modelo.

Fuente de datos

Selección de las covariables

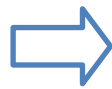
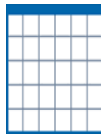
Transformación de las covariables

Ejemplo:



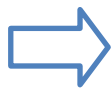
BURO  
(RCC)

Deudas  
de TC  
últimos  
6  
meses



$Deuda_{t-6}, Deuda_{t-5}, Deuda_{t-4},$   
 $Deuda_{t-3}, Deuda_{t-2}, Deuda_{t-1}$

Línea  
de  
TC



$LineaTC_{t-1}$

Calificac  
ión de  
Riesgo  
últimos 3  
meses



$Calif_{t-3}, Calif_{t-2}, Calif_{t-1}$

$$x_1 = \frac{Linea_{t-1}}{Deuda_{max}}$$

Ratio

$$x_2 = \frac{Deuda_{t-1} - Deuda_{t-6}}{Deuda_{t-6}}$$

Varianza

$$x_3 = Deuda_{promedio}$$

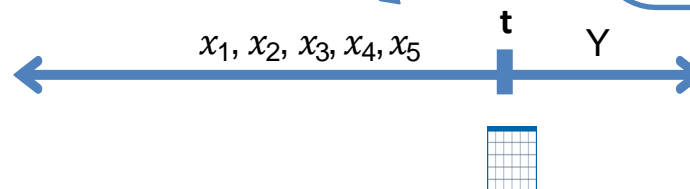
Tendencia

$$x_4 = Calif_{max}$$

Indicador

$$x_5 = VecesCPPaMas_{sum}$$

Frecuencia



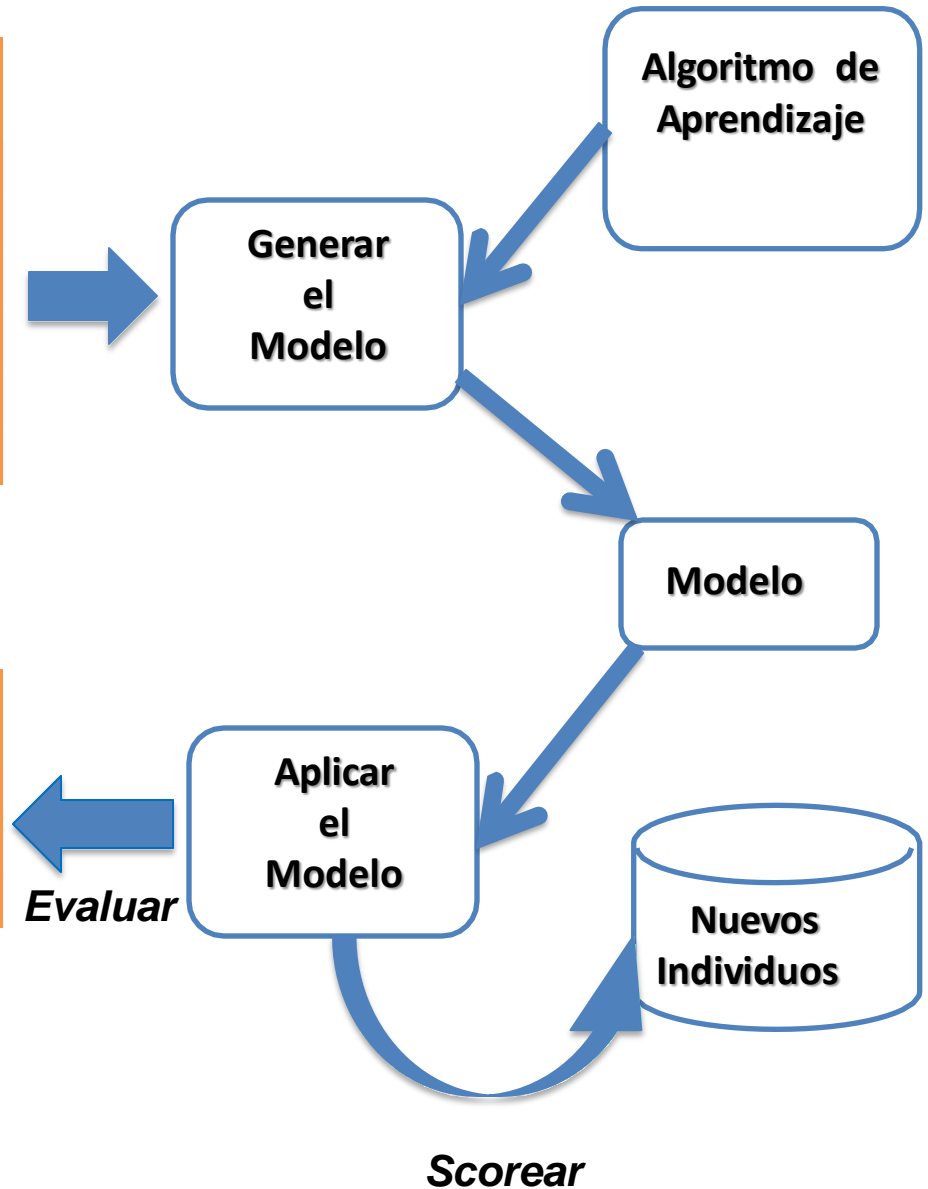
# Modelo general de Implementación de Modelos

ID	REEMBOLSO	ESTADO CIVIL	INGRESOS ANUALES	FRAUDE
1	SI	SOLTERO	S/ 1,000	NO
2	SI	CASADO	S/ 5,000	NO
3	NO	CASADO	S/ 3,500	SI
4	SI	VIUDO	S/ 4,500	NO
5	NO	SOLTERO	S/ 2,000	NO
6	NO	SOLTERO	S/ 1,500	SI

**Tabla de Aprendizaje**

ID	REEMBOLSO	ESTADO CIVIL	INGRESOS ANUALES	FRAUDE
7	SI	SOLTERO	S/ 4,000	NO
8	SI	CASADO	S/ 5,500	NO
9	NO	CASADO	S/ 6,500	SI

**Tabla de Testing**



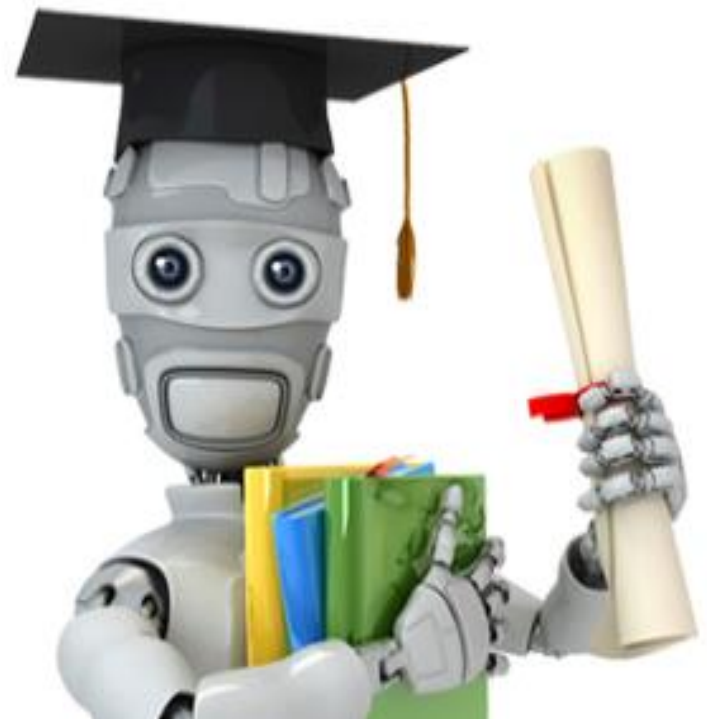
# METODOLOGÍA DE DESARROLLO DE ALGORITMOS PREDICTIVOS

- Definición y diseño del modelamiento a realizar.
- Análisis exploratorio de datos:
  - ✓ Detección de Outliers.
  - ✓ Imputación de valores perdidos.
  - ✓ Transformaciones.
  - ✓ Recodificaciones.
- Balanceo de datos.
- Selección de variables. (Met. Estadísticas vs ML).
- Modelamiento.
- Validación .
- Implementación.



# Implementación de Algoritmos de Machine Learning

30 %



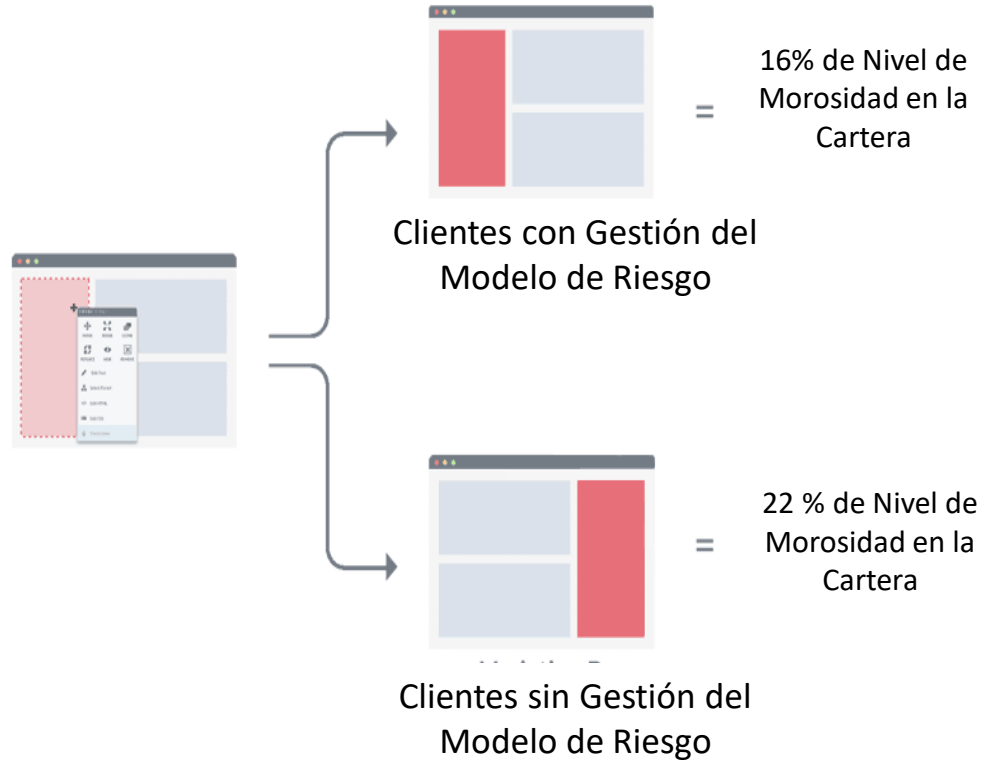
# Implementación de Modelos : Generación de Grupos de Ejecución

Probabilidad	Nº Clientes	Nº Sucesos VD	% Sucesos VD/ Nº Clientes
0,9	10 000	100	36%
0,8	10 000	60	22%
0,7	10 000	40	14%
0,6	10 000	33	12%
0,5	10 000	20	7%
0,4	10 000	10	4%
0,3	10 000	5	2%
0,2	10 000	5	2%
0,1	10 000	3	1%
0	10 000	3	1%
Total	100 000	279	100%

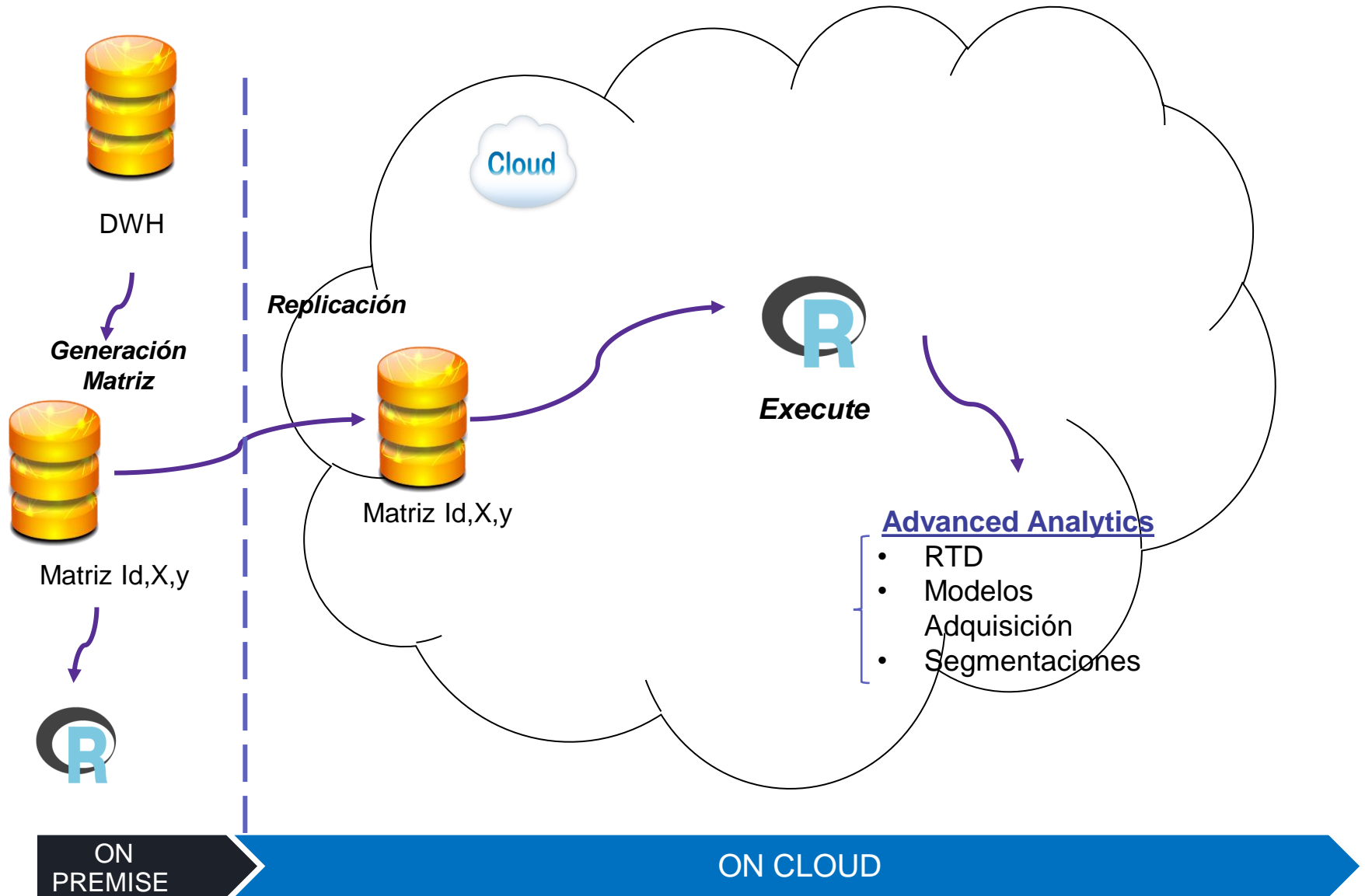


GRUPOS DE EJECUCIÓN	Nº CLIENTES	% SUCESOS ACUMULADOS	EFFECTIVIDAD	LIFT
<b>RECOMENDADO</b>	<b>30 000</b>	<b>200</b>	<b>0,67%</b>	<b>2,39</b>
MEDIO	30 000	63	0,21%	0,75
BAJO	40 000	16	0,04%	0,14
<b>TOTAL</b>	<b>100 000</b>	<b>279</b>	<b>0,28%</b>	

# Implementación de Modelos : A/B Testing



# Implementación de Modelos : Infraestructura



# ¡Gracias!

**André Omar Chávez Panduro**  
**UNMSM**

MSc in Data Science Candidate  
Promotion “Erwin Kraenau Espinal”  
**Universidad Ricard Palma**