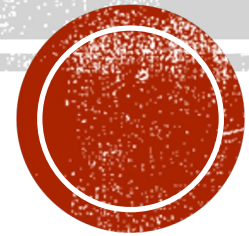


# PROYECTO APRENDIZAJE SUPERVISADO



“You can’t manage what you don’t measure”,  
Tom De Marco, 1982

“We’re drowning in data but starving for knowledge”,  
John Naisbitt, 1982

# DESCRIPCIÓN

En este curso corto (4 sesiones), utilizando un enfoque práctico, se repasarán los conceptos básicos de la analítica de datos, en especial los propios del aprendizaje supervisado (regresión y clasificación).

# DESCRIPCIÓN

## Aprendizaje Supervisado

- Tareas de regresión
- Tareas de clasificación

# UNIDADES

---

## **Unidad 1**

Aprendizaje supervisado  
Regresión

## **Unidad 2**

Aprendizaje supervisado  
Clasificación





# METODOLOGÍA

## Sesión de clase

- Repaso de los conceptos básicos a utilizar dentro del proyecto.
- Presentación y discusión del proyecto
- Revisión del conjunto de datos a utilizar en la tarea.
- Trabajo grupal (breakout rooms)



# ANDRÉS A. ARISTIZÁBAL P.

## **Formación**

- Ingeniero de Sistemas y Computación de la Universidad Javeriana Cali, 2006
- Doctorado en Informática, École Polytechnique de París, Francia, 2012

## **Experiencia académica**

- Investigador en el grupo Avispa de la Universidad Javeriana Cali, 2006 - 2009
- Investigador Posdoctoral en el grupo de Lenguajes de Programación de la Universidad de Wrocław, Polonia 2014 - 2015
- Profesor hora cátedra de la Universidad Javeriana Cali, 2015
- Profesor hora cátedra de la Universidad Icesi Cali, 2016 – 2017
- Desde 2017, profesor tiempo completo Universidad Icesi, Facultad de Ingeniería



# AGENDA

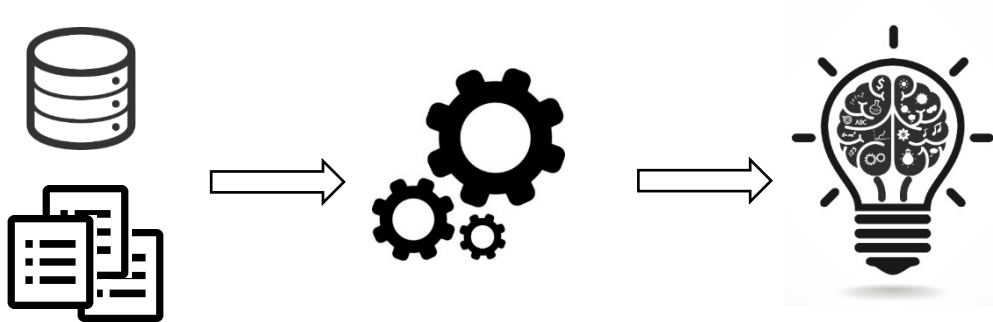
- Aprendizaje automático
- Aprendizaje supervisado
- Tareas de regresión
- Modelos de regresión
- Proyecto de regresión



# APRENDIZAJE AUTOMÁTICO

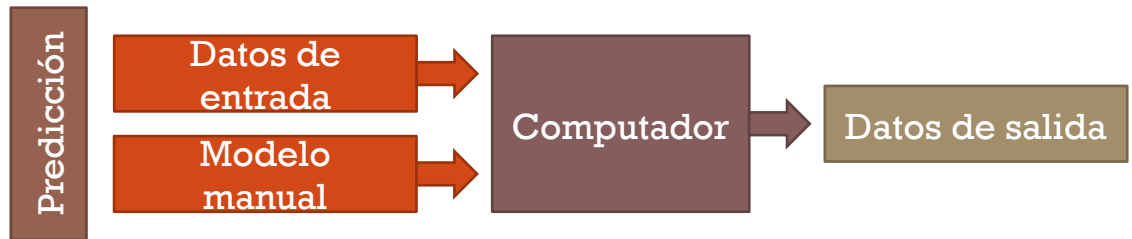
- Definición:

El aprendizaje automático es la ciencia que permite a los computadores aprender, sin ser explícitamente programados<sup>1</sup>

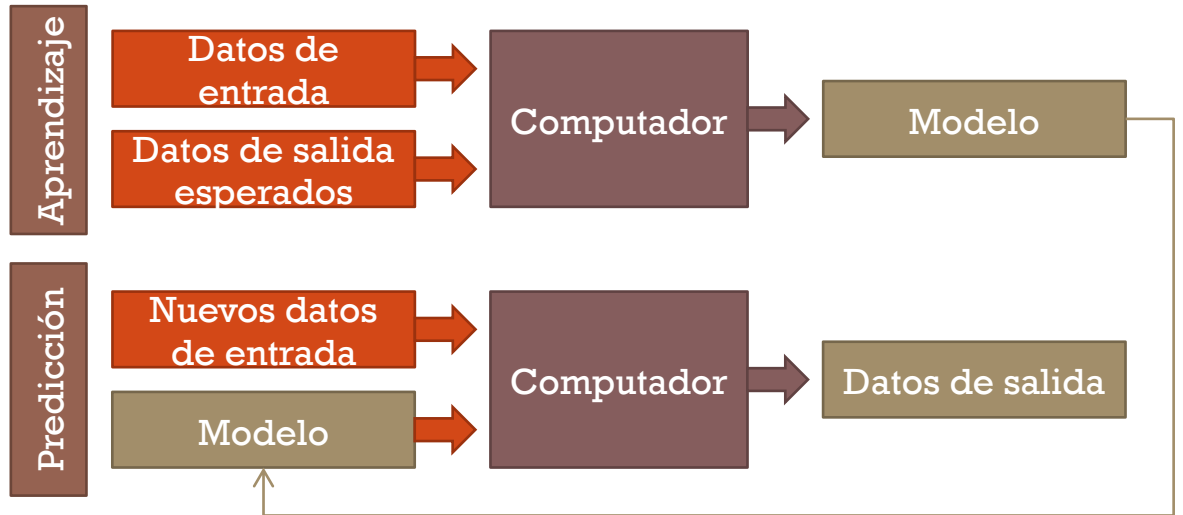


1. Andrew Ng, Stanford University, 2014

## Modelo tradicional



## Ciencia de datos





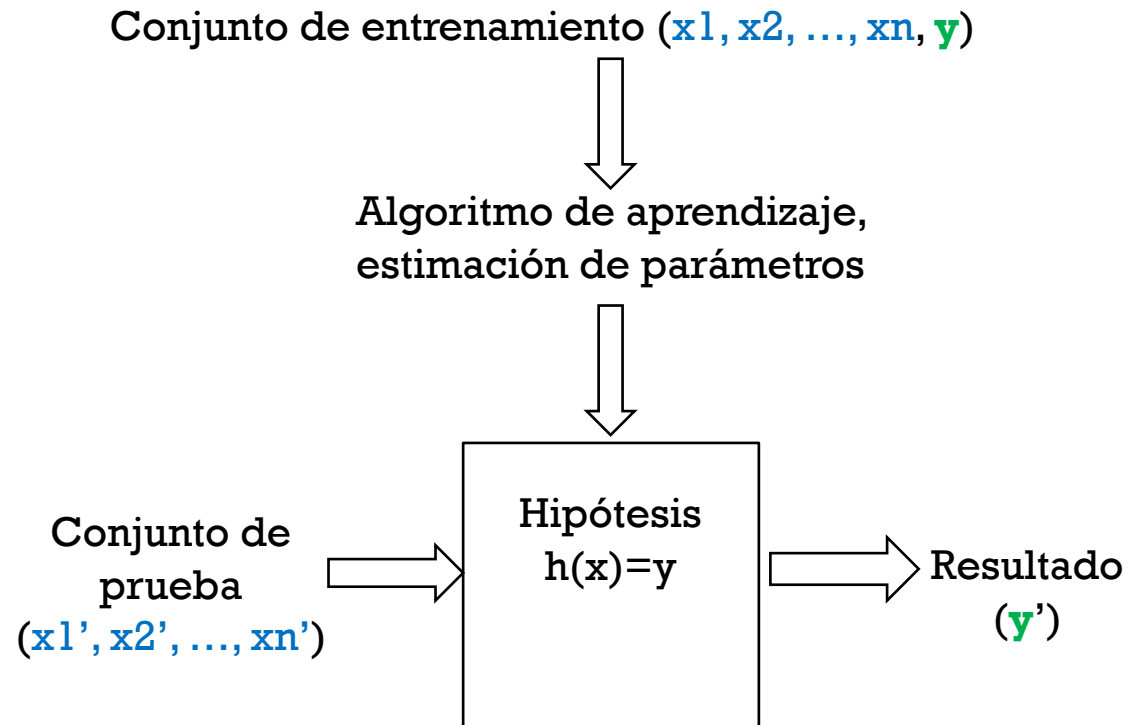
# APRENDIZAJE SUPERVISADO

- Aprender a partir de un “experto”
- Datos de entrenamiento **etiquetados** con una clase o valor:

$(x_1, x_2, \dots, x_n, y)$

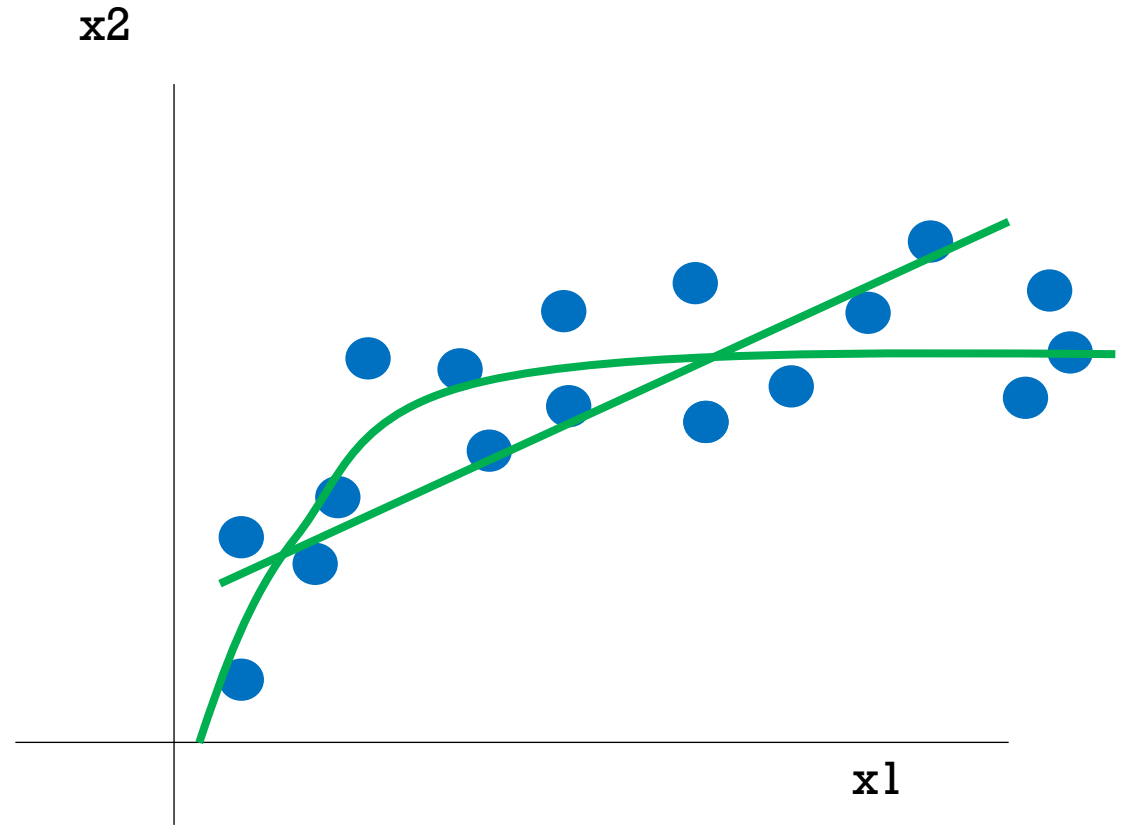
Predictores, variables de entrada (independientes)      Respuesta, variable de salida (dependiente)

- **Meta:** predecir una clase o valor



# TAREAS DE REGRESIÓN

- Encontrar modelos,  $f$ , que permitan **predecir valores continuos**:
  - KNN
  - Regresión lineal
  - Regresión polinómica
  - Árboles de regresión
  - ...
- Valores **numéricos** de la variable o función objetivo
- **Baseline**: medida de evaluación dada por un modelo que predice una medida de tendencia central (e.g. el promedio) (→ “modelo **nulo**”)



# MODELOS DE REGRESIÓN

- **Regresión lineal:** busca una relación lineal entre los atributos predictores ( $x_i$ ) y el atributo objetivo ( $Y$ )

$$Y = h_{\Theta}(X) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n + \varepsilon$$

Los parámetros  $\theta_i$  son estimados teniendo como objetivo la minimización de los residuos o diferencias cuadradas entre las predicciones ( $\hat{Y}$ ) y los valores reales ( $Y$ ):

$$\operatorname{argmin}_{\Theta} \sum_1^m (y - \hat{y})^2$$

- **Regresión Lasso:** regularización L1 (suma de los valores absolutos de los coeficientes multiplicados por un lambda)

$$\operatorname{argmin}_{\Theta} \sum_1^m (y - \hat{y})^2 + \lambda \sum_1^m |\theta_i|$$

- **Regresión Ridge:** regularización L2 (suma de los valores cuadrados de los coeficientes multiplicados por un lambda)

$$\operatorname{argmin}_{\Theta} \sum_1^m (y - \hat{y})^2 + \lambda \sum_1^m \theta_i^2$$

- **Regresión polinomial:** busca una relación polinomial entre los atributos predictores y la variable objetivo ( $Y$ )

$$Y = h_{\Theta}(X) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \dots + \theta_n x^n + \varepsilon$$



# MODELOS DE REGRESIÓN

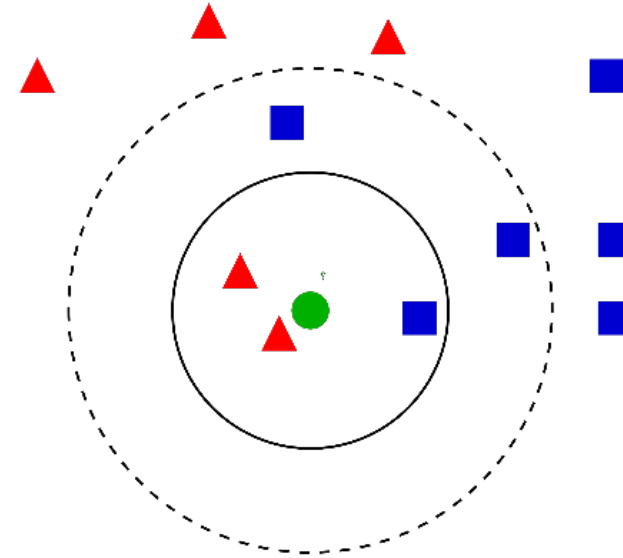
- **KNN**
- **Árboles de decisión**
- **Random Forest**
- **Boosting**



# MODELOS DE REGRESIÓN

## KNN:

- Algoritmo de aprendizaje supervisado para clasificación y regresión.
- Asigna la clase o valor agregado de las instancias conocidas que se encuentran mas cerca de la instancia a predecir.



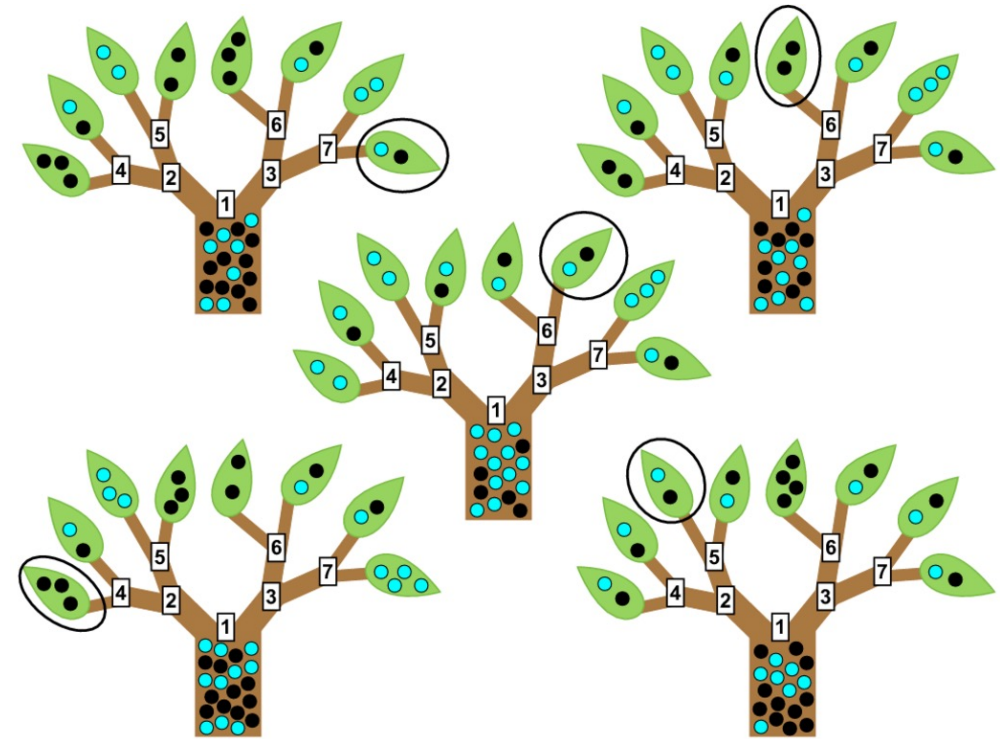




# MODELOS DE REGRESIÓN

## Random Forest:

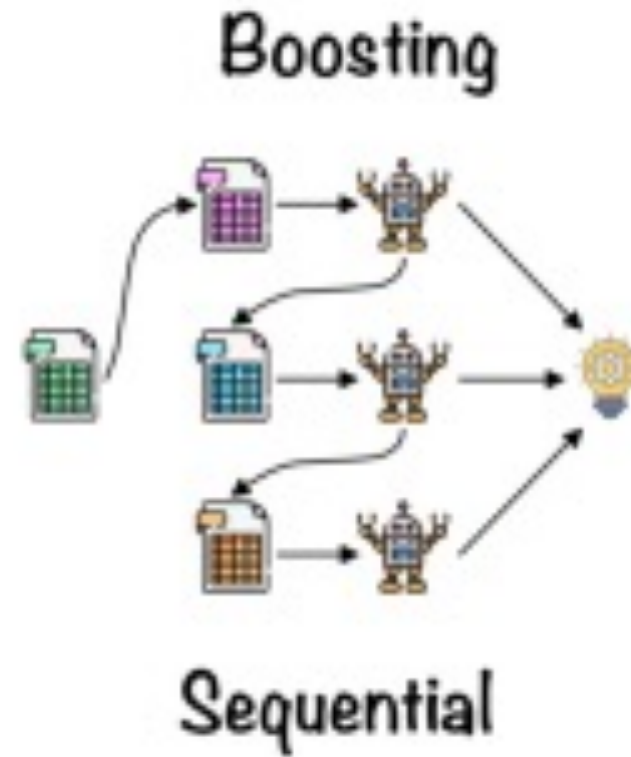
- Algoritmo de aprendizaje supervisado para clasificación y regresión.
- Basado en el concepto de ensambles, el cual es un proceso de combinación de múltiples regresores o clasificadores para resolver un problema complejo y mejorar el rendimiento del modelo.



# MODELOS DE REGRESIÓN

## Boosting:

- Algoritmo de aprendizaje supervisado para clasificación y regresión.
- Mejora la exactitud y rendimiento de los modelos de aprendizaje automático convirtiendo múltiples modelos débiles en uno mucho más potente.



# Proyecto de regresión

- Exploración y visualización
- Ingeniería de características
- Limpieza de datos
  - Imputación de variables
- Protocolos de evaluación
- Métricas de evaluación
- Modelo base
- Uso de pickle para guardar un modelo y sus métricas
- Modelo de regression lineal
- Selección de características
- Modelos Lasso y Ridge
- Modelos polinomiales
- Modelo de K vecinos más cercanos
- Modelo de árboles de decisión
- Modelo de Random forest
- Modelos de Boosting
- Comparación de modelos

