# MODELOS DE REGRESIÓN

Manuel Valdivia



### Data Analytics en Banca – Modelamiento de los datos

### 1. Regresión Lineal

La regresión lineal simple consiste en generar un modelo de regresión (ecuación de una recta) que permita explicar la relación lineal que existe entre dos variables. A la variable dependiente o respuesta se le identifica como Y, a la variable predictora o independiente como X.

El modelo de regresión lineal simple se describe de acuerdo a la ecuación:

$$Y = \beta_0 + \beta_1 X_1 + \epsilon$$

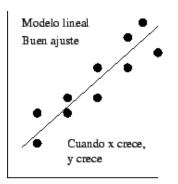
### Siendo:

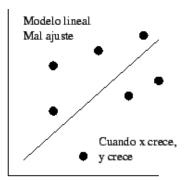
- la ordenada en el origen  $eta_0$
- la pendiente
- el error aleatorio.

Este último representa la diferencia entre el valor ajustado por la recta y el valor real.

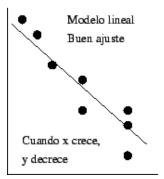
### Data Analytics en Banca – Modelamiento de los datos

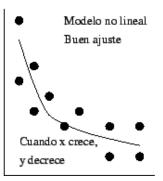
### 1. Regresión Lineal













### Data Analytics en Banca – Modelamiento de los datos

### 1. Regresión Lineal

### Condiciones:

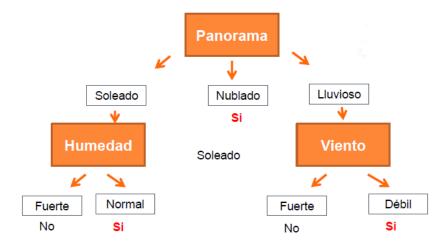
- Linealidad: La relación entre ambas variables debe ser lineal.
- Distribución Normal de los residuos: Los residuos se tiene que distribuir de forma normal, con media igual a 0.
- ❖ Varianza de residuos constante (homocedasticidad): La varianza de los residuos ha de ser aproximadamente constante a lo largo del eje X.
- Valores atípicos y de alta influencia: Hay que estudiar con detenimiento los valores atípicos o extremos ya que pueden generar una falsa correlación que realmente no existe, u ocultar una existente.
- Independencia, Autocorrelación: Las observaciones deben ser independientes unas de otras.

### Data Analytics en Banca – Modelamiento de los datos

### 2. Arboles de Clasificación, Arboles de Regresión

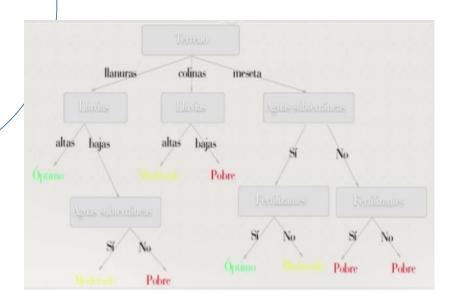
Un Árbol de decisión es un modelo de predicción que divide el espacio de los predictores agrupando observaciones con valores similares para la variable respuesta o dependiente. Dado un conjunto de datos se fabrican diagramas de construcciones lógicas, muy similares a los sistemas de predicción basados en reglas, que sirven para representar y categorizar una serie de condiciones que ocurren de forma sucesiva

### **Ejemplo:** la pregunta a responder es Juego Tennis?



### Data Analytics en Banca – Modelamiento de los datos

2. Arboles de Clasificación, Arboles de Regresión



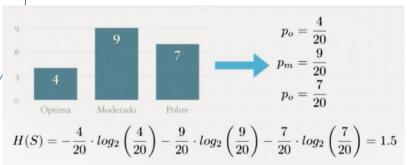
- Cosecha = Óptima
  - Terreno = Llanura y Lluvias = Altas
  - 0 Terreno = Meseta y Aguas Subterráneas = Sí y Fertilizantes = Sí
- Cosecha = Moderada
  - Terreno = Llanura y Lluvias = Bajas y Aguas Subterráneas=/Sí
  - 0 Terreno = Colinas y Lluvias = Altas
- o Terreno = Meseta y Aguas Subterráneas = Sí y Fertilizantes = No

Llúvias	Óptima	Moderada	Pobre	Terreno	Óptima	Moderada	Pobre
Altas	3 (43%)	4 (57%)	0	Colinas	0	3 (50%)	3 (50%
	1 (17%)	2 (33%)	3 (50%)		2 (33%)	1 (17%)	3 (50%
	0	2 (29%)	5 (71%)		2 (25%)	3 (37.5%)	3 (37.5%
Fertilizante	Óptima	Moderada	Pobre	Agua Sub.	Óptima	Moderada	Pobre
Fertilizante	Óptima	Moderada	Pobre	Agua Sub.	Óptima	Moderada	Pobre
Fertilizante No	Óptima 0	Moderada 6 (67%)	Pobre 3 (33%)	Agua Sub.	Óptima 2 (18%)	Moderada 6 (55%)	Pobro

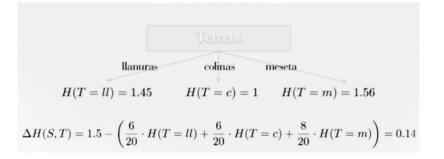
### Data Analytics en Banca – Modelamiento de los datos

2. Arboles de Clasificación, Arboles de Regresión

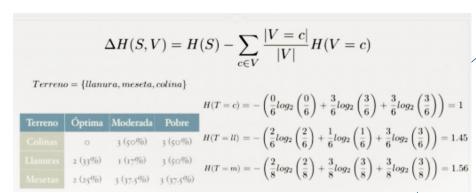
### Entropía



### Ganancia de Información



### Ganancia de Información



### Ganancia de Información

Parámetro	Ganancia de Información
	0.14
	0.42
	0.36
	0.16

### Data Analytics en Banca – Modelamiento de los datos

2. Arboles de Clasificación, Arboles de Regresión

### Índice de Gini

Para cada	variable	tenemos
Categoría	Cı	C2
	2	6
	8	4
	10	IO

### Reducción de la Varianza

# Para cada variable tenemosCategoríaCIC2TotalI268Convertimoso8412Si = ITotalIOIO2O

### Reducción de la Varianza

Para cada	variab	le llevan	nos a cabo	$mean(root) = \frac{8 \cdot 1 + 12 \cdot 0}{20} = 0.4$ $var(root) = \frac{8 \cdot (1 - 0.4)^2 + 12 \cdot (0 - 0.4)^2}{20} = 0.24$
Categoría	Cı	C <sub>2</sub>	Total	$mean(C_1) = \frac{2 \cdot 1 + 8 \cdot 0}{10} = 0.2$
1	2	6	8	$Var(C_1) = \frac{2 \cdot (1 - 0.2)^2 + 8 \cdot (0 - 0.2)^2}{10} = 0.16$
	8	4	12	$mean(C_2) = \frac{6 \cdot 1 + 4 \cdot 0}{10} = 0.6$
	10	IO	20	$Var(C_2) = \frac{6 \cdot (1 - 0.6)^2 + 4 \cdot (0 - 0.6)^2}{10} = 0.24$
				$Var.Ponderada = \frac{10}{20} \cdot 0.16 + \frac{10}{20} \cdot 0.24 = 0.20$

Data Analytics en Banca – Modelamiento de los datos

2. Arboles de Clasificación, Arboles de Regresión

### Algoritmos:

- CART
- CHAID
- ID3
- C50

### Data Analytics en Banca – Modelamiento de los datos

### Algoritmo CART

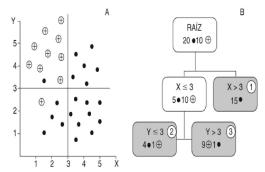
### **Notas:**

- Árboles de clasificación: predicen categorías de objetos.
- Árboles de regresión: predicen valores continuos.
- Partición binaria recursiva.
- En cada iteración se selecciona la variable predictiva y el punto de separación que mejor reduzcan la 'impureza'.

### Características:

- Uno de los métodos de aprendizaje supervisado no paramétrico más utilizado.
- Realizar sucesivas divisiones binarias en un conjunto de datos guiado por un criterio.
- Para cada nodo selecciona a la variable independiente que proporciona el mejor desempeño en el criterio para particionar los datos.

### Idea Intuitiva:



### Data Analytics en Banca – Modelamiento de los datos

### **Algoritmo CART**

### Criterios de Partición:

- Objetivo de una partición: Incrementar la homogeneidad (en términos de clase) de los subconjuntos resultantes. Que sean más puros que el conjunto originario.
- Que sean más puros que el conjunto originario. Existen criterios de impureza tales como : Medida de Entropia, Indice de Gini.

### Índice de Gini:

$$G(A_i) = \sum_{j=1}^{M_i} p(A_{ij})G(C|A_{ij})$$

$$G(C|A_{ij}) = -\sum_{k=1}^{J} p(C_k|A_{ij})p(\neg C_k|A_{ij}) =$$

$$= 1 - \sum_{k=1}^{J} p^2(C_k|A_{ij})$$

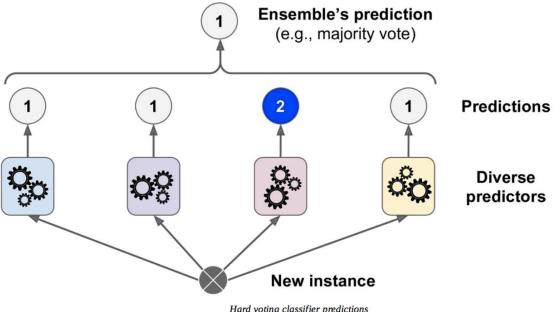
### Criterio de Parada:

- Un nodo se declarará terminal si el nodo es puro.
- Un nodo se declarará terminal si cualquier otra subdivisión no da una mejora mayor que la obtenida en el nodo padre.

- Ai es el atributo para ramificar el árbol.
- Mi es el número de valores diferentes del atributo Ai.
- p(Aij) es la probabilidad de que Ai tome su j-ésimo valor (1 <= j <= Mi).</li>
- p(Ck | Aij) es la probabilidad de que un ejemplo pertenezca a la clase Ck cuando su atributo Ai toma su j-ésimo valor.

Data Analytics en Banca – Modelamiento de los datos

3. Ensemble Voting

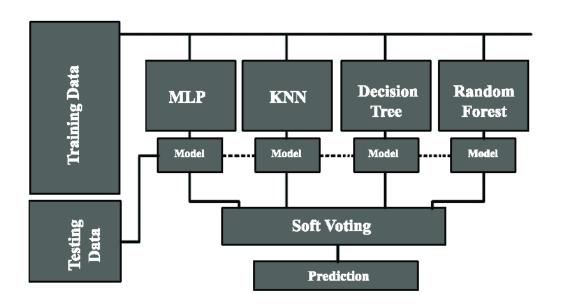


Hard voting classifier predictions

Es una técnica que se puede utilizar para mejorar el rendimiento del modelo, idealmente logrando un mejor rendimiento que cualquier modelo único utilizado en el conjunto.

### Data Analytics en Banca – Modelamiento de los datos

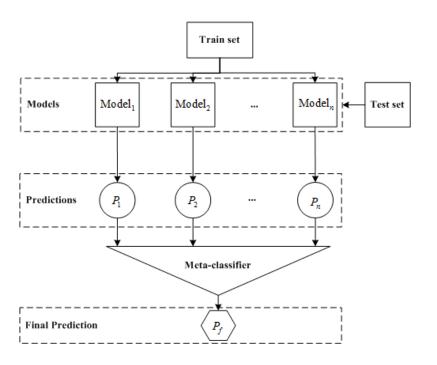
3. Ensemble Voting



En el caso de la regresión, esto implica calcular el promedio de las predicciones de los modelos. En el caso de la clasificación, se suman las predicciones para cada etiqueta y se predice la etiqueta con el voto mayoritario.

### Data Analytics en Banca – Modelamiento de los datos

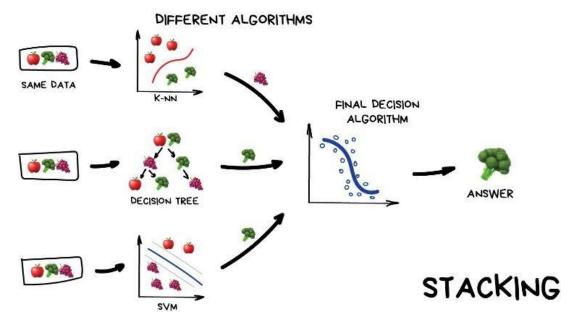
4. Ensemble Stacking



Utiliza un algoritmo de metaaprendizaje para aprender cómo combinar mejor las predicciones de dos o más algoritmos básicos de aprendizaje automático.

### Data Analytics en Banca – Modelamiento de los datos

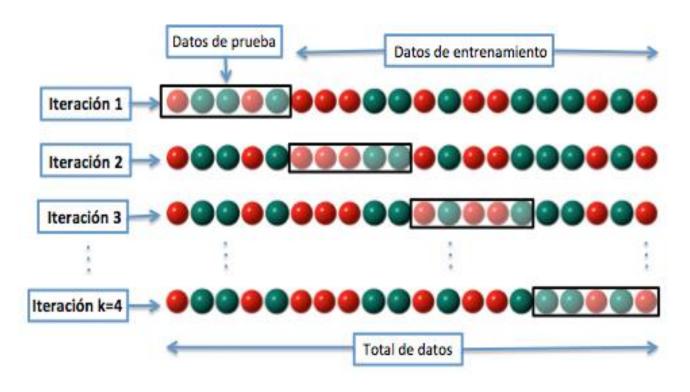
### 4. Ensemble Stacking



El beneficio del apilamiento es que puede aprovechar las capacidades de una variedad de modelos con buen desempeño en una tarea de clasificación o regresión y hacer predicciones que tienen un mejor desempeño que cualquier modelo individual en el conjunto.

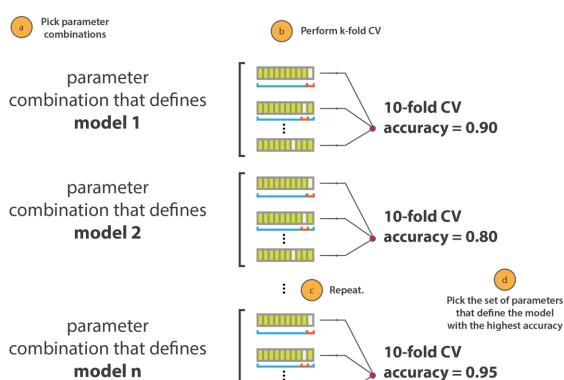
Data Analytics en Banca – Modelamiento de los datos

5. Cross Validation



### Data Analytics en Banca – Modelamiento de los datos

### 5. Cross Validation



# **GRACIAS**

