



Universidad **Ricardo Palma**

RECTORADO

PROGRAMA DE ESPECIALIZACIÓN EN CIENCIA DE DATOS

Formamos seres humanos para una cultura de paz

PROGRAMA DE ESPECIALIZACIÓN DATA SCIENCE NIVEL I



R + Python

MÓDULO I



**Modelos de Regresión Logística,
Naive Bayes y K-NN**



A nuestro recordado Maestro

Dr. Erwin Kraenau Espinal, Presidente de la Comisión de Creación de la Maestría en Ciencia de los Datos



PROGRAMA DE ESPECIALIZACIÓN EN DATA SCIENCE NIVEL I

« Para poder **seguir** a veces hay que **empezar de nuevo**»



Agenda

- Introducción a los Problemas de Clasificación o Aprendizaje Supervisado.
- Regresión Logística Binaria.
- Regresión Logística Binaria Penalizada.
- K – Vecinos más cercanos.
- Elección del K – óptimo.
- Clasificación Bayesiana : Naive Bayes.



CLASIFICACIÓN: DEFINICIÓN

- Dada una colección de registros (Conjunto de Entrenamiento) cada registro contiene un conjunto de variables (atributos) denominado x , con una variable (atributo) adicional que es la clase denominada y .
- El objetivo de la **clasificación** es encontrar un modelo (una función) para predecir la clase a la que pertenecería cada registro, esta asignación una clase se debe hacer con la mayor precisión posible.
- Un conjunto de prueba (tabla de testing) se utiliza para determinar la precisión del modelo. Por lo general, el conjunto de datos dado se divide en dos conjuntos al azar de el de entrenamiento y el de prueba.



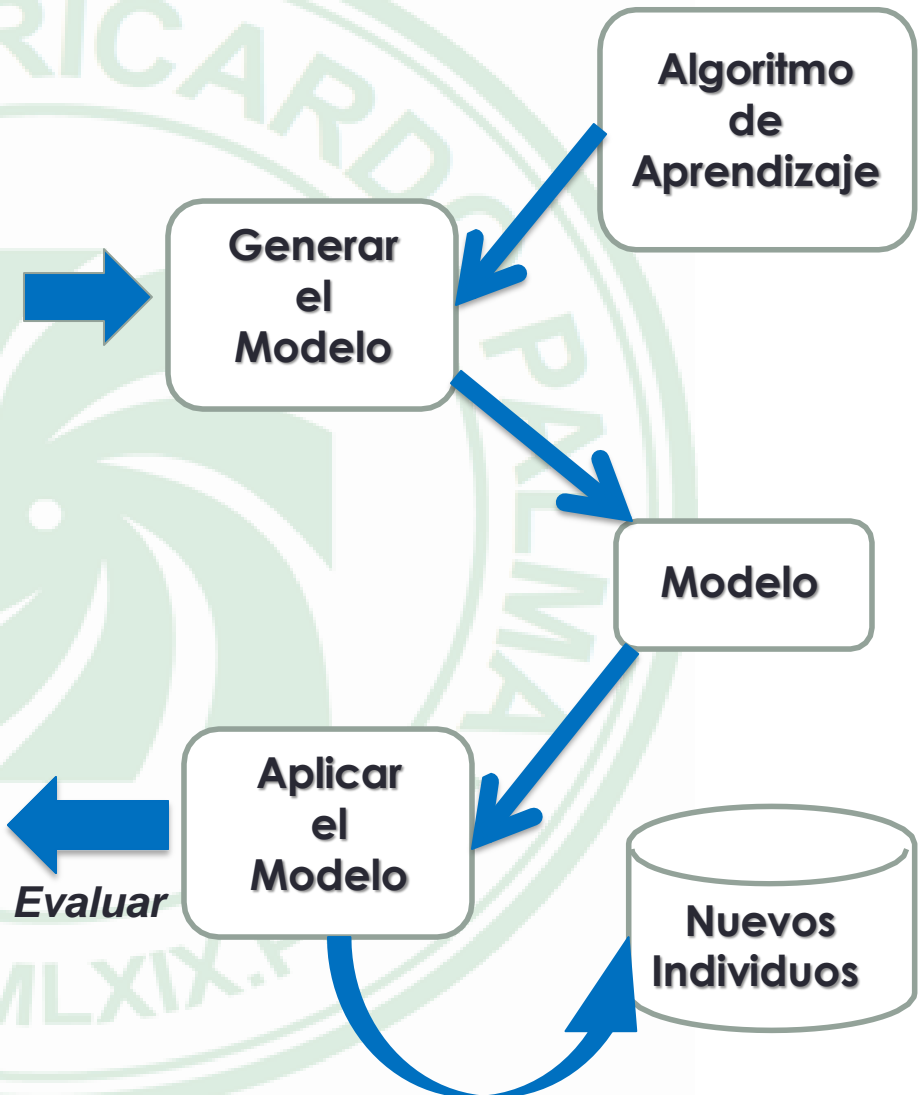
MODELO GENERAL DE LOS MÉTODOS DE CLASIFICACIÓN

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

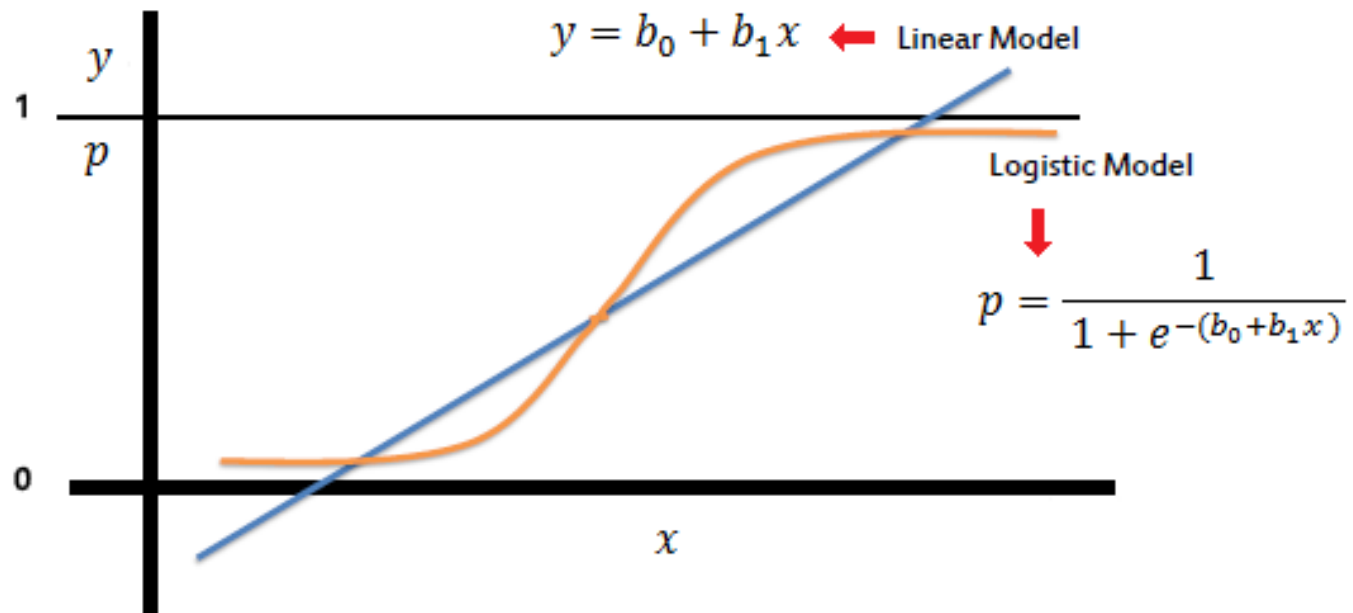


DEFINICIÓN DE CLASIFICACIÓN

- Dada una base de datos $D = \{t_1, t_2, \dots, t_n\}$ de tuplas o registros (individuos) y un conjunto de clases $C = \{C_1, C_2, \dots, C_m\}$, el **problema de la clasificación** es encontrar una función $f: D \rightarrow C$ tal que cada t_i es asignada una clase C_j .
- $f: D \rightarrow C$ podría ser una Red Neuronal, un Árbol de Decisión, un modelo basado en Análisis Discriminante, o una Red Bayesiana.



Regresión Logística



DEFINICIÓN

- Es un modelo predictivo supervisado.
- La regresión logística es un modelo de elección discreta en el que la variable dependiente es cualitativa.
- Es flexible en cuanto a la naturaleza de las variables explicativas, pues éstas pueden ser de escala y categóricas.



MODELO DE REGRESIÓN LOGÍSTICA DICOTÓMICO

- Para este modelo se considera que la variable respuesta, es una variable dicotómica que toma dos valores.
- Para estos modelos dicotómicos, las dos categorías deben de ser mutuamente excluyentes.
- La variable respuesta se puede expresar de la siguiente forma:

$$Y_i = \begin{cases} 1, \text{Prob}(Y_i = 1) = P_i \\ 0, \text{Prob}(Y_i = 0) = 1 - P_i \end{cases}$$



MODELO DE REGRESIÓN LOGÍSTICA DICOTÓMICO

Ejemplo:

La variable Morosidad toma los siguientes valores:

- “1” si el cliente es moroso.
- “0” si el cliente es no moroso.

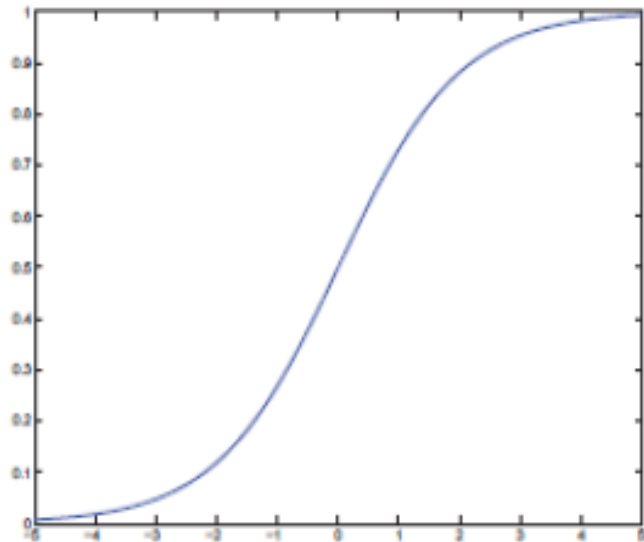
¿Es dicotómica?

¿Es cualitativa?

¿Es mutuamente excluyente?



MODELO DE REGRESIÓN LOGÍSTICA DICOTÓMICO



Se basa en la función logística:

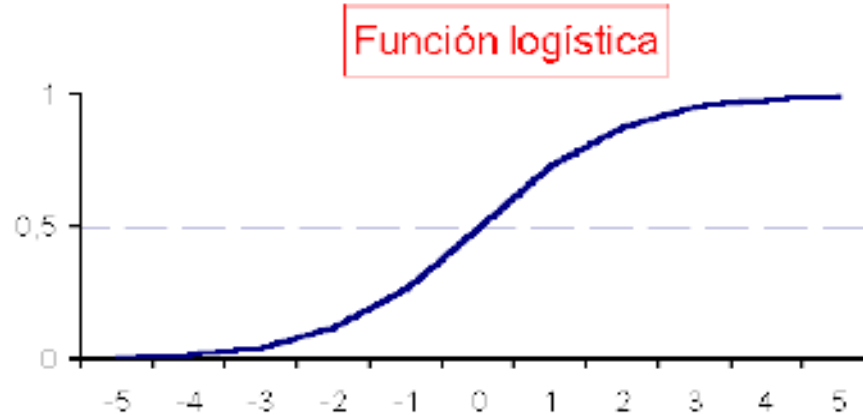
$$f(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + \frac{1}{e^z}} = \frac{e^z}{1 + e^z}$$

Está acotada entre 0 y 1:

$$\lim_{z \rightarrow -\infty} f(z) = 0, \quad \lim_{z \rightarrow \infty} f(z) = 1,$$



MODELO DE REGRESIÓN LOGÍSTICA DICOTÓMICO



La representación matemática del modelo es la siguiente:

$$z_i = \log \frac{P_i}{1 - P_i} = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

z_i : Variable dependiente del modelo: “Moroso” y “No Moroso”

p_i : Probabilidad de que el cliente sea “Moroso”

β_i : Coeficientes del modelo (parámetros a estimar)

x_i : Variables explicativas del modelo



MODELO DE REGRESIÓN LOGÍSTICA DICOTÓMICO

Odds Ratio

Es la razón entre la probabilidad de que se produzca un suceso y la probabilidad de que no se produzca ese suceso.

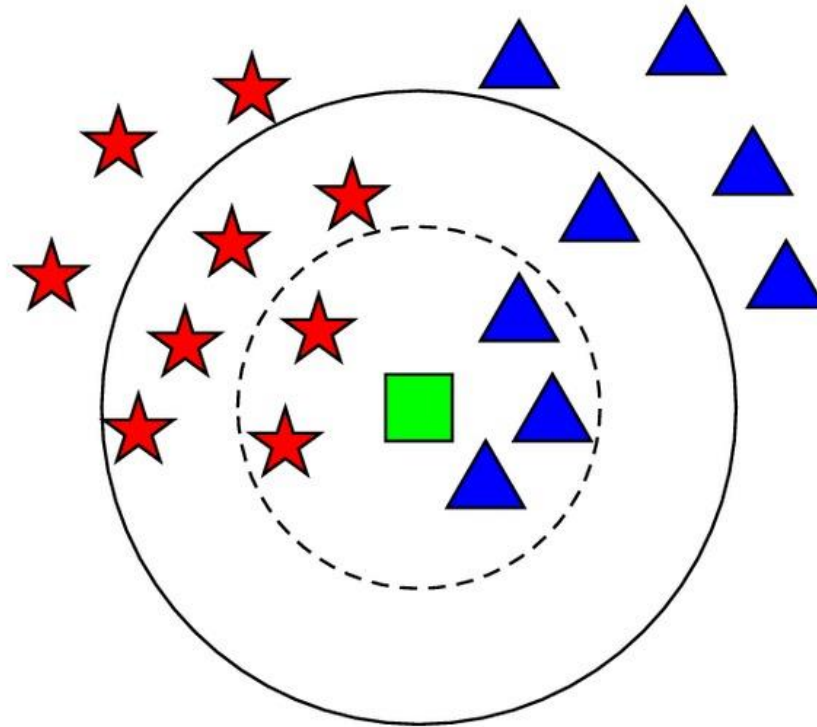
$$z_i = \log \frac{P_i}{1 - P_i} = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$



Indica cuánto más probable es ser un cliente “Moroso” que “No Moroso”

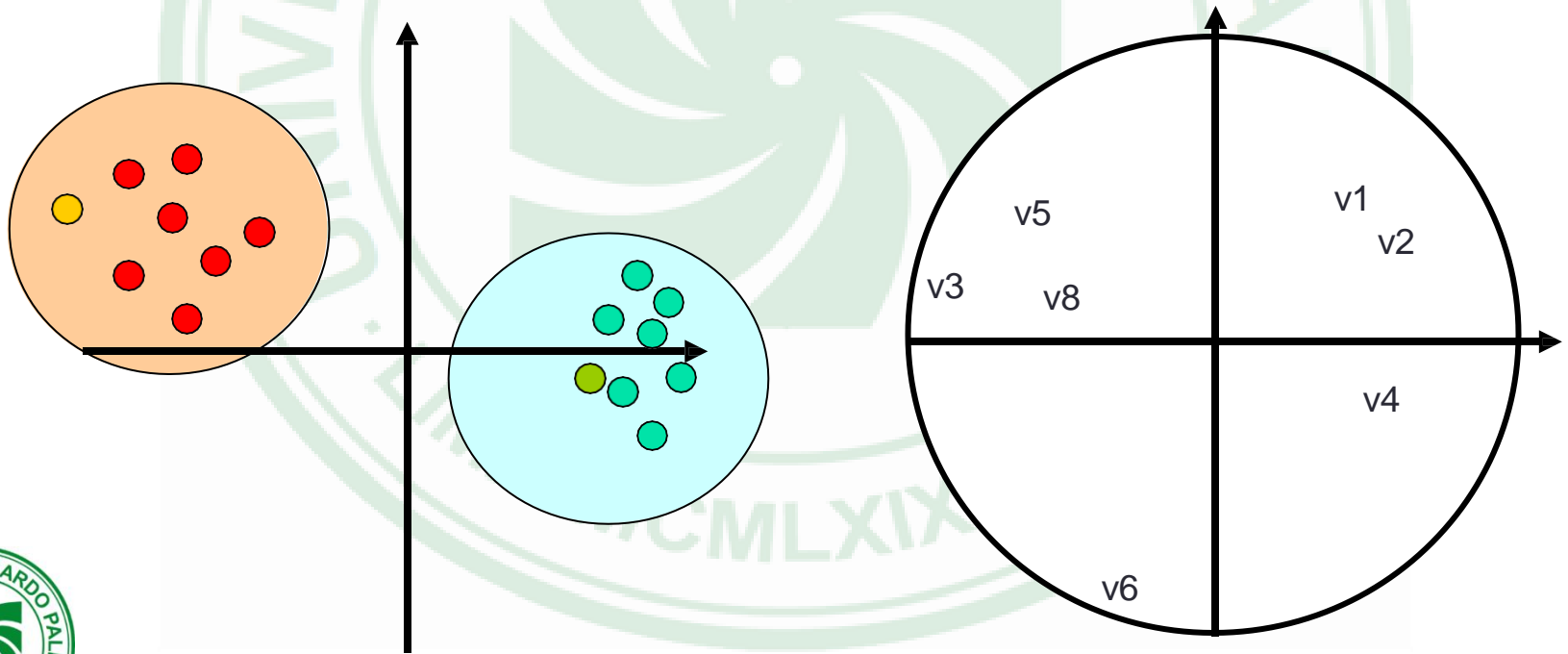


K – Vecinos más cercanos



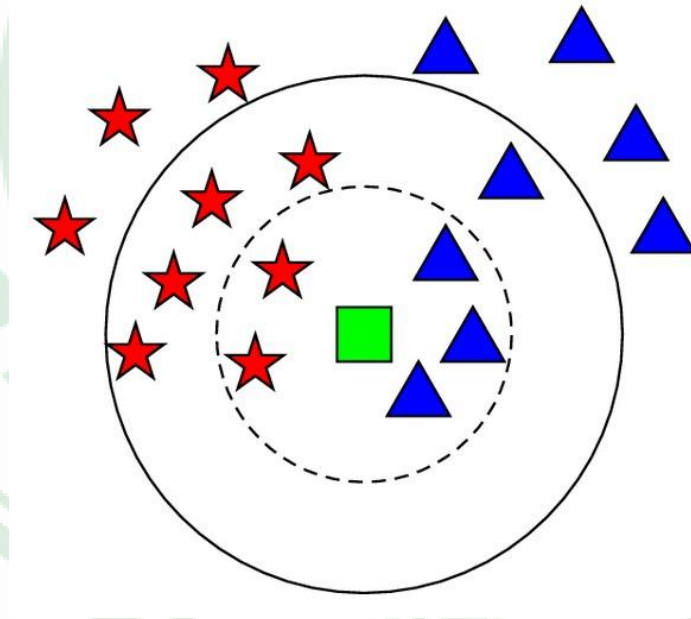
CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS

- Análisis de vecino más próximo es un método de clasificación de casos basado en su similaridad con otros casos. En aprendizaje de máquinas, se ha desarrollado como una forma de reconocer patrones de datos sin requerir una coincidencia exacta con patrones o casos almacenados.



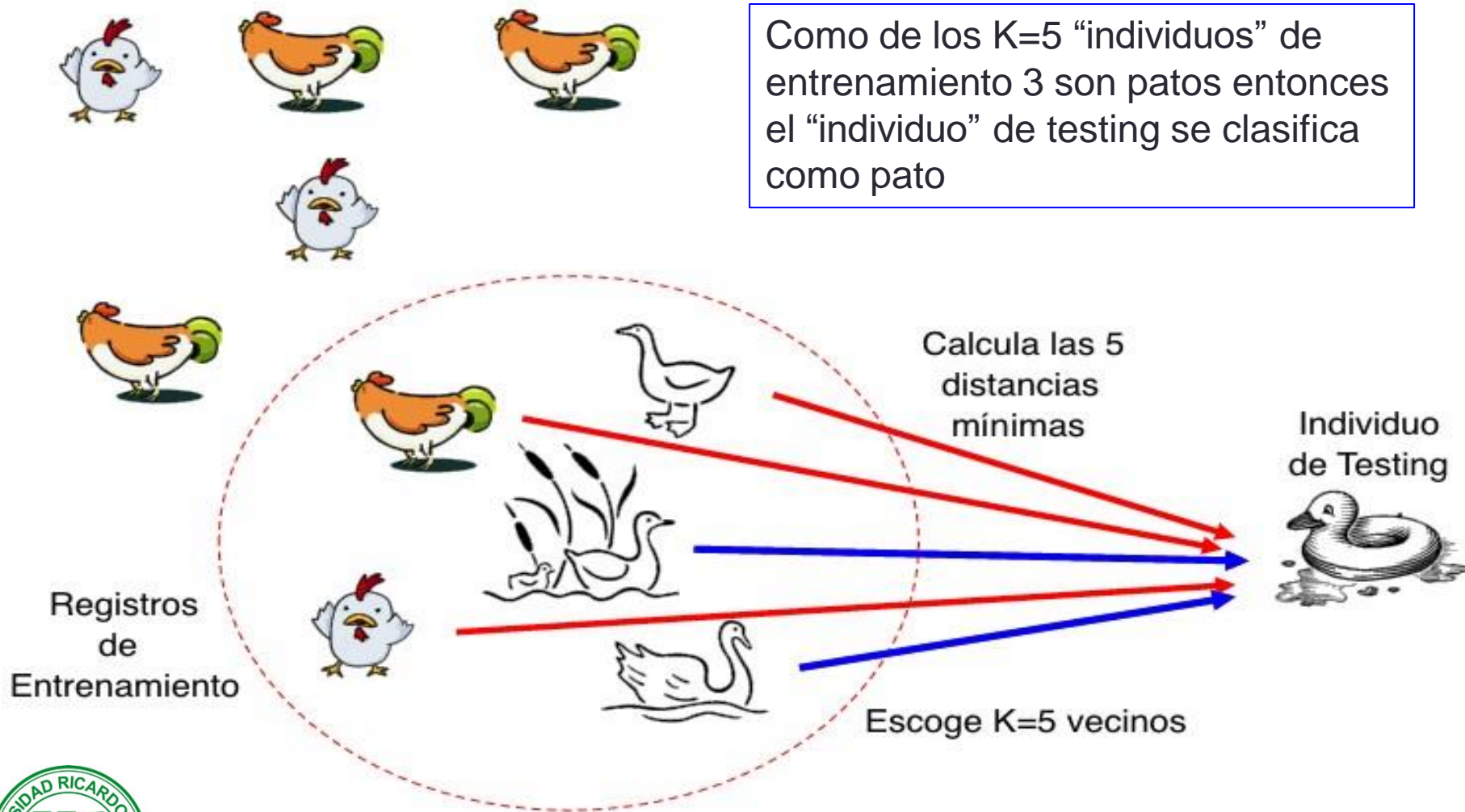
CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS

- Los casos similares están cercanos entre sí y los casos no similares están distantes entre sí. Además, la distancia entre dos casos es una medida de sus diferencias.

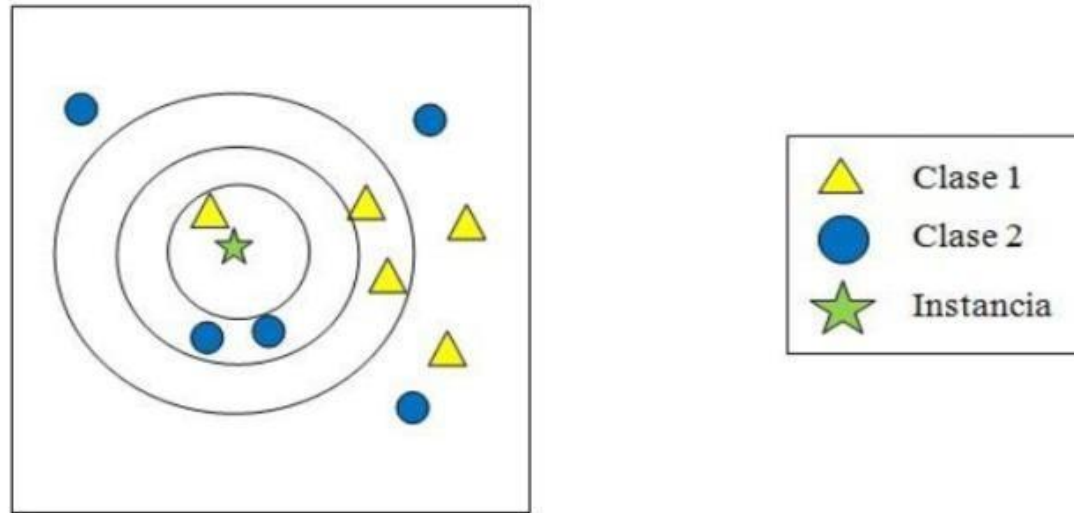


CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : IDEA INTUITIVA

Como de los $K=5$ "individuos" de entrenamiento 3 son patos entonces el "individuo" de testing se clasifica como pato



CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS



Para $K=1$ (círculo más pequeño), la clase de la nueva instancia sería la Clase 1, ya que es la clase de su vecino más cercano, mientras que para $K=3$ la clase de la nueva instancia sería la Clase 2 pues habrían dos vecinos de la Clase 2 y solo 1 de la Clase 1

CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ALGORITMO

COMIENZO

Entrada: $D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_N, c_N)\}$

$\mathbf{x} = (x_1, \dots, x_n)$ nuevo caso a clasificar

PARA todo objeto ya clasificado (x_i, c_i)

calcular $d_i = d(\mathbf{x}_i, \mathbf{x})$

Ordenar $d_i (i = 1, \dots, N)$ en orden ascendente

Quedarnos con los K casos $D_{\mathbf{x}}^K$ ya clasificados más cercanos a \mathbf{x}

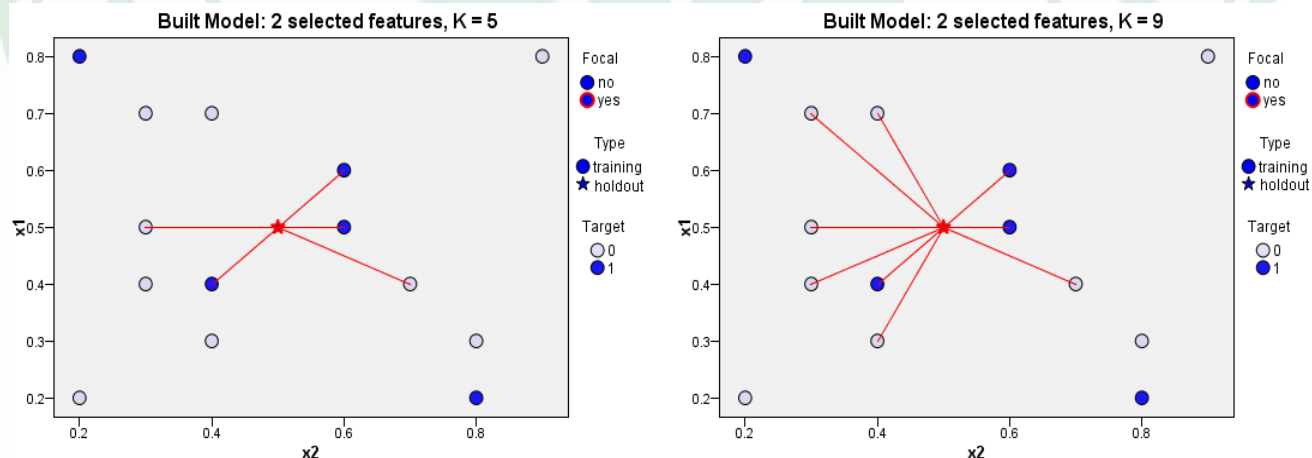
Asignar a \mathbf{x} la clase más frecuente en $D_{\mathbf{x}}^K$

FIN

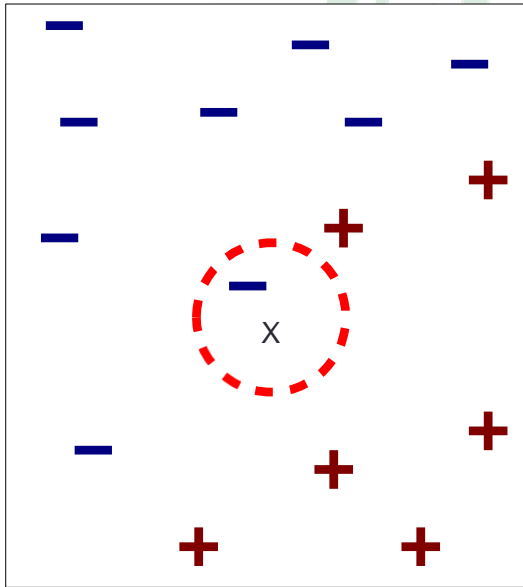


CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ALGORITMO

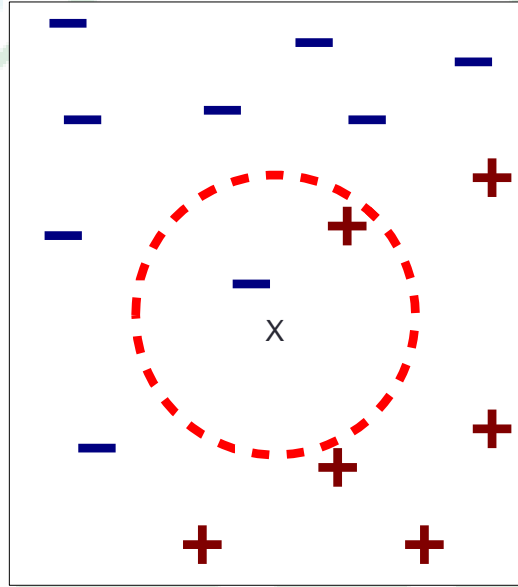
- Los casos muy cercanos a otros se denominan “vecinos”. Cuando se presenta un nuevo caso (reserva), se calcula su distancia desde cada caso del modelo. Las clasificaciones de la mayoría de casos similares (los vecinos más próximos) se anotan y el nuevo caso se coloca en la categoría que contiene el mayor número de vecinos más próximos.
- Puede especificar el número de vecinos más próximos que se van a examinar; este valor se denomina k .



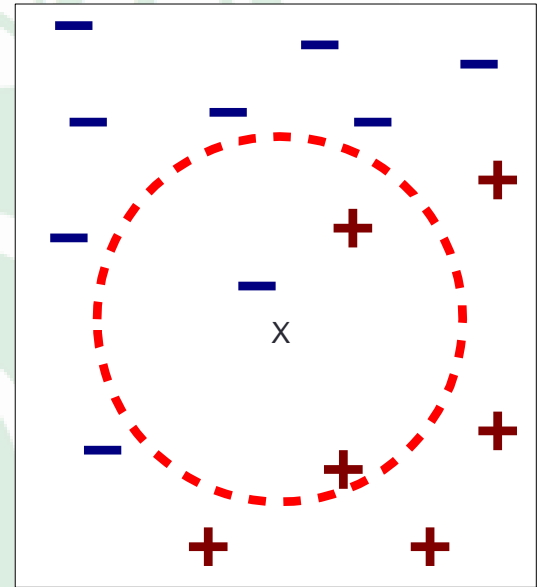
CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ¿ELECCIÓN DEL K ÓPTIMO?



(a) 1 - Vecino más cercano.



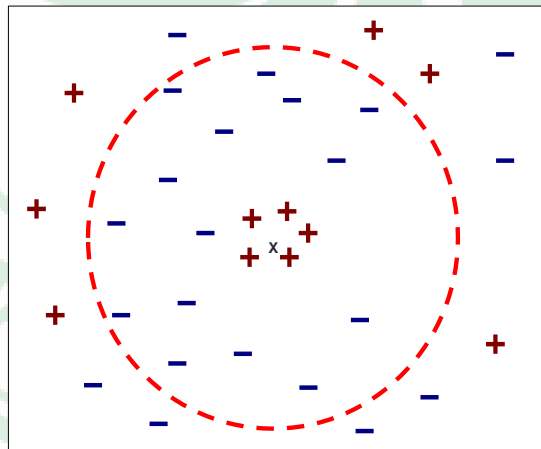
(b) 2 - Vecinos más cercanos.



(c) 3- Vecinos más cercanos.

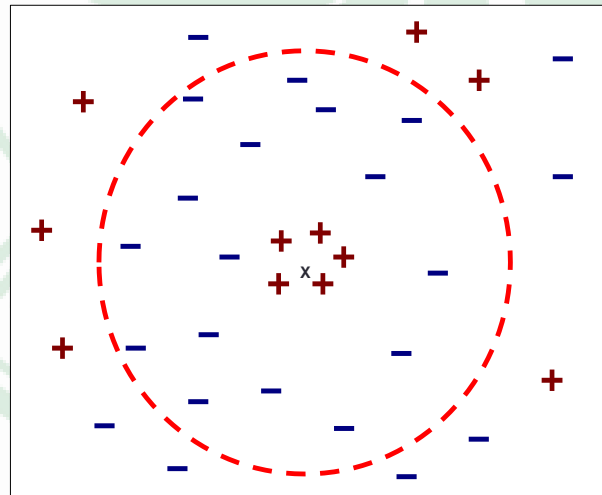
CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ¿ELECCIÓN DEL K ÓPTIMO?

- Escogiendo el valor de K:
 - ✓ Si K es muy pequeño el modelo será muy sensitivo a puntos que son atípicos o que son ruido (datos corruptos)
 - ✓ Si K es muy grande, el modelo tiende a asignar siempre a la clase más grande.



CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ¿ELECCIÓN DEL K ÓPTIMO?

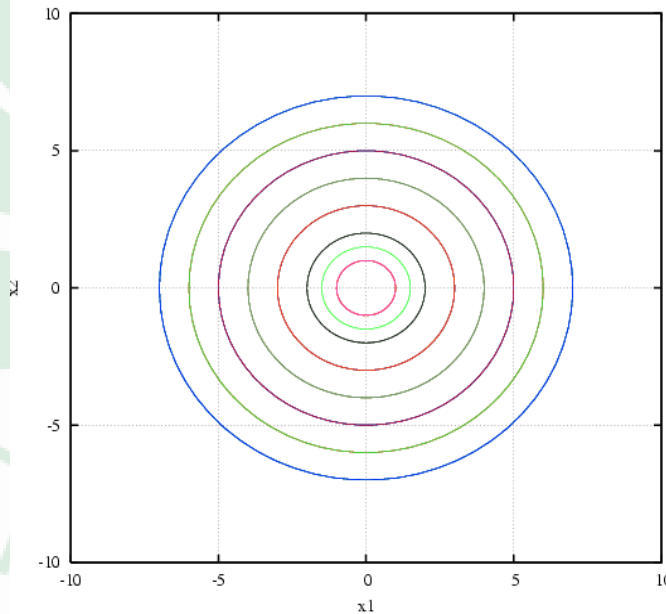
- Escogiendo el valor de K:
 - ✓ Mediante la Tabla de Aprendizaje el modelo escogerá el valor de K que mejor clasificación logre en esta tabla, es decir, prueba con $K=1$, $K=2$,
 - ✓ Esto puede ser muy caro computacionalmente.



CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ELECCIÓN DE LA DISTANCIA

$$d(A, B) \equiv \sqrt{\sum_{i=1}^n (A_i - B_i)^2} = \sqrt{(A - B)^T (A - B)}$$

**Distancia
Euclídea**



CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : DISTANCIA EUCLÍDEA

$$d(A, B) \equiv \sqrt{\sum_{i=1}^n (A_i - B_i)^2} = \sqrt{(A - B)^T (A - B)}$$

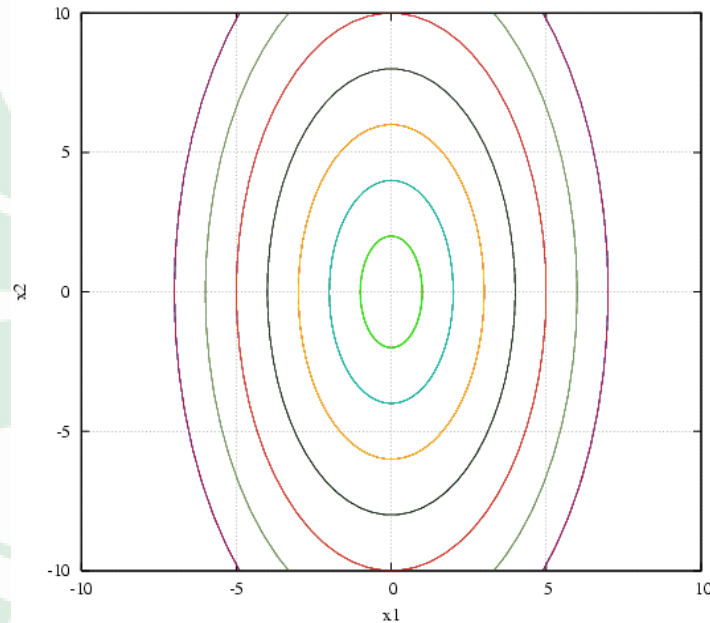
- Algunas consideraciones importantes para el uso de ésta distancia, es que no está acotada, es sensible a cambios de escalas y considera las T variables estocásticamente independientes.



CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ELECCIÓN DE LA DISTANCIA

$$d(A, B) \equiv \sqrt{(A - B)^T M^T M (A - B)}$$

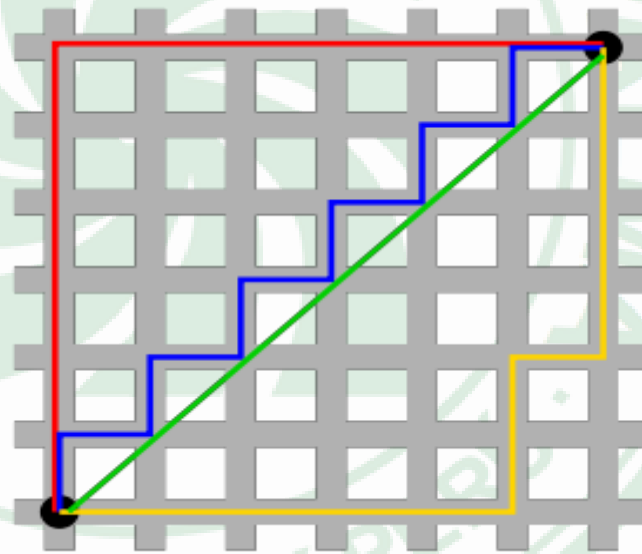
**Distancia
Euclídea
Ponderada**



CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ELECCIÓN DE MÉTRICA DE DISTANCIA

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

Distancia
Manhattan
(city-block)



CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : DISTANCIA DE MANHATTAN

- La distancia de Manhattan o métrica city-block o ciudad, calculada como la suma de las diferencias absolutas entre unidades para cada variable, es un caso particular de la distancia de **Minkowski**. Es menos sensible a valores muy grandes o aberrantes, ya que es función de diferencias absolutas en lugar de diferencias al cuadrado, adicionalmente cada variable puede ser estandarizada por su rango.

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

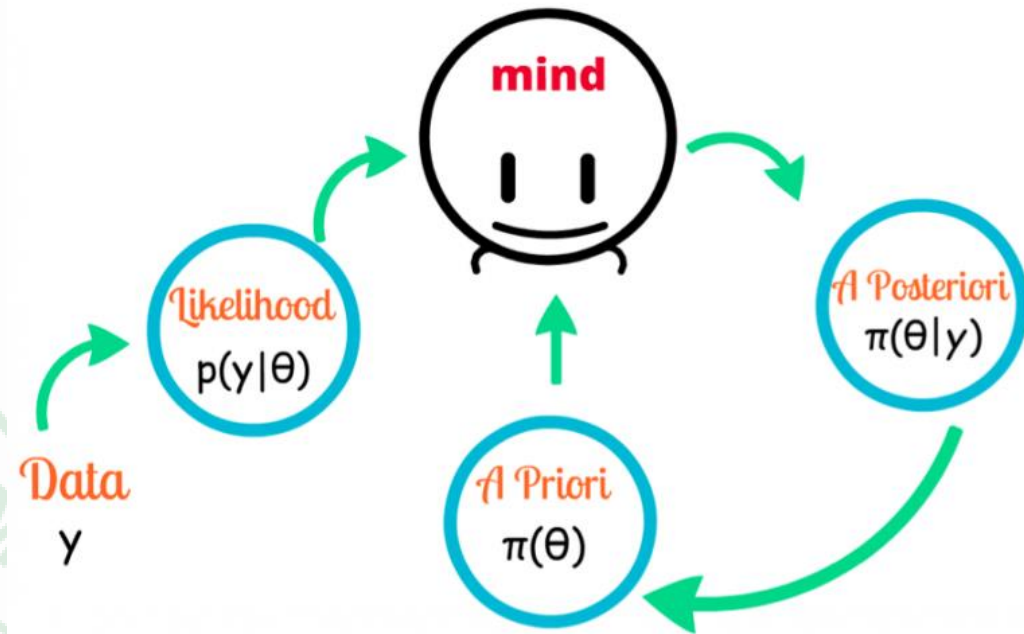


CLASIFICACIÓN MEDIANTE K-VECINOS MÁS CERCANOS : ELECCIÓN DE LA DISTANCIA

Distancias y disimilaridades		Métrica	Euclídea
Euclídea	$\sqrt{\sum_{i=1}^T (x_{ii} - x_{ji})^2}$	Si	Si
Manhattan	$\sum_{i=1}^T x_{ii} - x_{ji} $	Si	No
Bray-Curtis	$\frac{\sum_{i=1}^T x_{ii} - x_{ji} }{\sum_{i=1}^T (x_{ii} + x_{ji})}$	Si	No
Canberra	$\sum_{i=1}^T \frac{ x_{ii} - x_{ji} }{(x_{ii} + x_{ji})}$	Si	No
Minkowski	$\sqrt[q]{\sum_{i=1}^T x_{ii} - x_{ji} ^q}$	Si	Si
Mahalanobis	$\sqrt{\sum_{i=1}^T \sum_{j=1}^T (x_{ii} - x_{ji}) \sigma_{ii}^{-1} (x_{ii} - x_{ji})}$	Si	Si



Modelos Bayesianos : Naive Bayes



INTRODUCCIÓN

- Estudió el problema de la determinación de la probabilidad de las causas a través de los efectos observados.
- El teorema que lleva su nombre se refiere a la probabilidad de un suceso condicionado por la ocurrencia de otro suceso.



Thomas Bayes



DEFINICIÓN

- Es un método importante no sólo porque ofrece un análisis cualitativo de las atributos y valores que pueden intervenir en el problema, sino porque da cuenta también de la importancia cuantitativa de esos atributos.
- En el aspecto cualitativo podemos representar cómo se relacionan esos atributos ya sea en una forma causal, o señalando simplemente de la correlación que existe entre esas variables (o atributos). Cuantitativamente (y ésta es la gran aportación de los métodos bayesianos).



DEFINICIÓN

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

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$



EJEMPLO PARA VARIABLES CUALITATIVAS

Ejemplo 1.(atributos discretos solamente)

X1	X2	X3	Y
0	0	1	0
0	1	0	0
1	1	0	0
0	0	1	1
1	1	1	1
0	0	1	1
1	1	0	1



EJEMPLO PARA VARIABLES CUANTITATIVAS

Clasificador Naïve Bayes (cont)

En este caso

$$P[X_j = a_j / C_i] = \frac{1}{s_j \sqrt{2\pi}} \exp\left[-\frac{(a_j - \bar{x}_j)^2}{2s_j^2}\right]$$

Donde, \bar{x}_j y s_j son la media y la varianza de los valores de la variable X_j en la clase C_i .

La libreria e1071 de R contiene una funcion **naiveBayes** que calcula el clasificador naïve Bayes, tanto para datos discretos como continuos.





¡Gracias!



Comunidad Data Science Perú



Comunidad Data Science Perú

PROGRAMA DE ESPECIALIZACIÓN EN DATA SCIENCE NIVEL I