

Histograma - Probabilidades

↳ Ejemplo correo spam.



$p(\text{Dear} | \text{N}) = 0.47$
 $p(\text{Friend} | \text{N}) = 0.29$
 $p(\text{Lunch} | \text{N}) = 0.18$
 $p(\text{Money} | \text{N}) = 0.06$

Now, because these histograms are taking up a lot of space, let's get rid of them, but keep the probabilities.



$p(\text{Dear} | \text{S}) = 0.29$
 $p(\text{Friend} | \text{S}) = 0.14$
 $p(\text{Lunch} | \text{S}) = 0.00$
 $p(\text{Money} | \text{S}) = 0.57$

Likelihoods / Probabilities

"Dear friend"

$$p(\text{N}) = 0.67$$

$$p(\text{N}) = \frac{8}{8+4} = 0.67$$

prior probability

0.09

Normal message

Normal message

$$p(\text{N}) \times p(\text{Dear} | \text{N}) \times p(\text{Friend} | \text{N})$$

$$p(\text{S}) \times p(\text{Dear} | \text{S}) \times p(\text{Friend} | \text{S})$$

0.01

2

Lunch Money Money Money Money

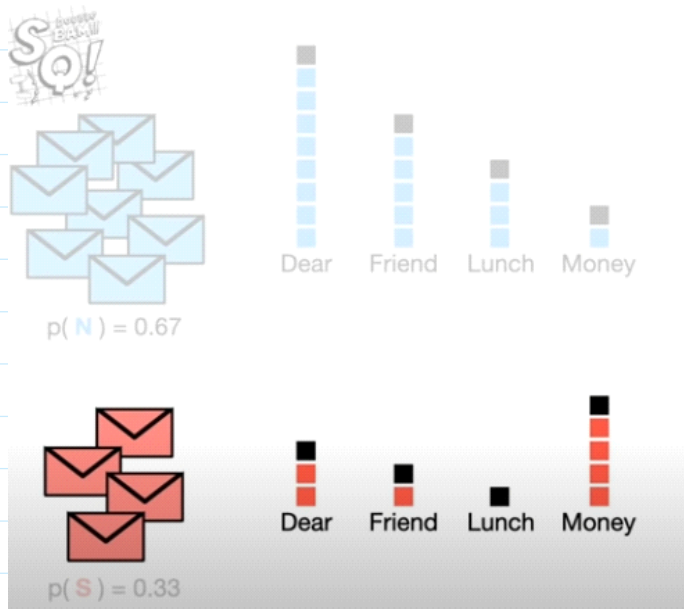
$$p(\text{N}) \times p(\text{Lunch} | \text{N}) \times p(\text{Money} | \text{N})^4 \approx 0.0002$$

$$p(\mathbf{N}) \times p(\mathbf{Lunch} | \mathbf{N}) \times p(\mathbf{Money} | \mathbf{N})^4 = 0.00002$$

$$p(\mathbf{S}) \times p(\mathbf{Lunch} | \mathbf{S}) \times p(\mathbf{Money} | \mathbf{S})^4 = 0 = 0$$

Para evitar esto

=> dumb ??



util //

$$p(\mathbf{N}) \times p(\mathbf{Lunch} | \mathbf{N}) \times p(\mathbf{Money} | \mathbf{N})^4 = 0.00001$$

$$p(\mathbf{S}) \times p(\mathbf{Lunch} | \mathbf{S}) \times p(\mathbf{Money} | \mathbf{S})^4 = 0.00122$$

↓ ¡SPAM!

Naive Bayes.

The thing that makes **Naive Bayes** so *naive* is that it treats all word orders the same.

Score for **Dear Friend** =

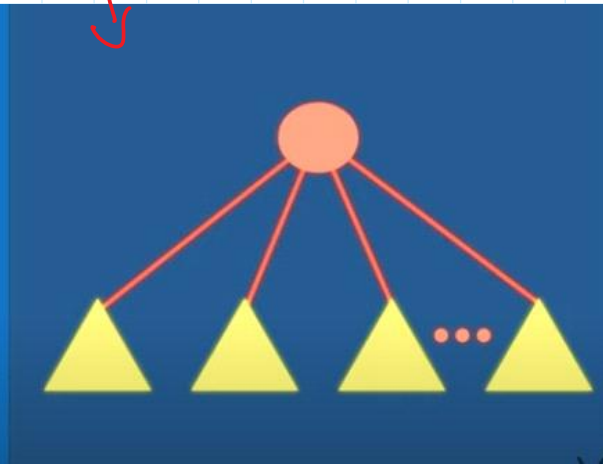
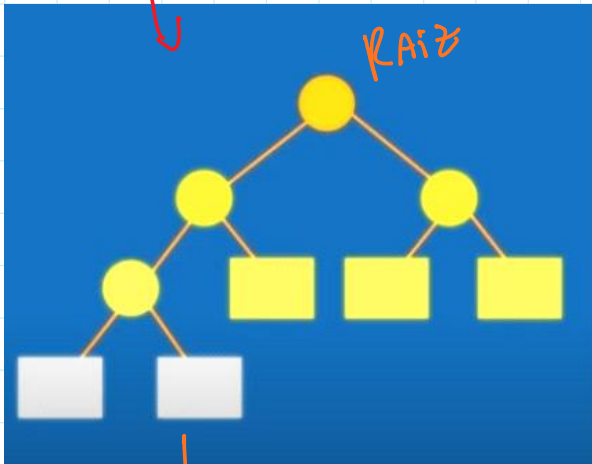
$$p(\mathbf{N}) \times p(\mathbf{Dear} | \mathbf{N}) \times p(\mathbf{Friend} | \mathbf{N}) = 0.08$$

Score for **Friend Dear** =

$$p(\mathbf{N}) \times p(\mathbf{Friend} | \mathbf{N}) \times p(\mathbf{Dear} | \mathbf{N}) = 0.08$$

Ignora las reglas de lenguaje

Decision Tree y Random Forest



Preguntas,
Arbol decision

Dataset / Etiquetas
y rasgos

				E
1	0	1	0	1
1	1	0	0	0
1	1	1	0	1

IN completas

Gini Impurity,

menor coeficiente
↳ Nodo raíz

Arbol de decision

Datos al Azar

**BOOTSTRAPPED
DATASET**

Arbol de decision
~ veces

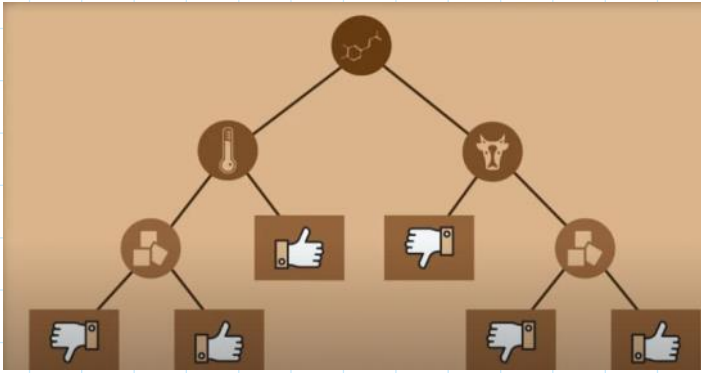
varios tambien

||

Dato clasificado
veces

(n)
100

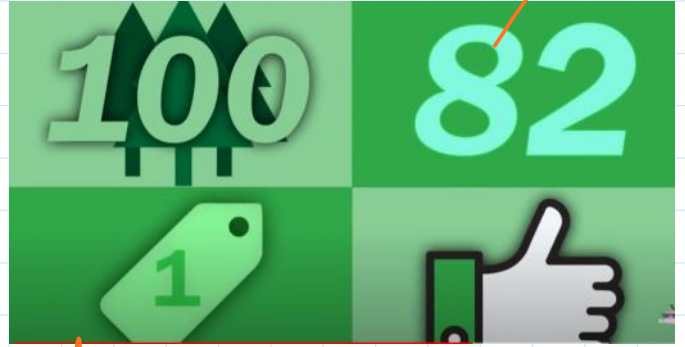
↳ Nodo raíz



↓
Clasificar Datos

↳ Numéricos y Booleanos

↳ Un clasificador



↳ Dependiendo de los n
árboles

↳ "Pena cruzada"

↳ Arbolito mejor
criterio

↳ Datos de prueba

↳ se pueden comparar los
árboles

↳ Varios clasificadores

↳ Difícil escoger