# Naive Bayes

* Classification
  $\hookrightarrow$ Multi-class

* $18^{th}$ Century $\longrightarrow$ Thomas Bayes $\Rightarrow$ Bayes Theorm
  Probability concepts

* "Naive" $\Rightarrow$ multicollinearity X

  Assumption

| Email | Discount $(X_1)$ | Mathematical Win $(X_2)$ | Spam/Not spam $(Y)$ |
|-------|------------------|--------------------------|---------------------|
| 1 | Y | Y | Spam |
| 2 | Y | N | Not spam |
| 3 | N | Y | Spam |
| 4 | Y | Y | Spam |
| 5 | N | N | Not spam |
| 6 | Y | Y | Spam |
| 7 | N | Y | Not spam |

Two phases $\longrightarrow$ Training , Prediction

Training
Phase    Step ① : Calculate probability of categories of
                  o/p colum

$$P = \frac{\text{No of obs of that category}}{\text{Total no of observation}}$$

(a) Spam

$$P(spam) = \frac{4}{7}$$

(b) Not spam

$$P(not\ spam) = \frac{3}{7}$$

Step ② : Calculate Conditional Probability / Likelihood

for spam

(i) $P(X_1 = \text{Yes} \mid \text{spam}) = \dfrac{X_1 = \text{Yes but belong to spam}}{\text{Total no of spam}}$

$$= \dfrac{3}{4}$$

(ii) $P(X_1 = \text{No} \mid \text{spam}) = \dfrac{1}{4}$

(iii) $P(X_2 = Yes \mid Spam) = \dfrac{4}{4}$

(iv) $P(X_2 = No \mid Spam) = 0$

for Not ~~Spam~~

(i) $P(X_1 = Yes \mid Not\ Spam) = \dfrac{1}{3}$

(ii) $P(X_1 = No \mid Not\ Spam) = \dfrac{2}{3}$

(iii) $P(X_2 = \text{Yes} \mid \text{Not spam}) = \dfrac{1}{3}$

(iv) $P(X_2 = \text{No} \mid \text{Not spam}) = \dfrac{2}{3}$

Training
Phase Completed

$\underline{\text{Prediction}}$
$\underline{\text{Phase}}$
$\begin{cases} X_1 = \text{Yes} \\ X_2 = \text{Yes} \end{cases}$ new data
$(\underline{\text{Bayes}} \ \underline{\text{Theorm}})$

$\underline{\text{Spam}}$

$P(\text{Spam} \mid X_1 = \text{Yes}, X_2 = \text{Yes}) = P(\text{spam}) \times P(X_1 = \text{Yes} \mid \text{spam})$

$\times P(X_2 = \text{Yes} \mid \text{Spam})$

$= \dfrac{4}{7} \times \dfrac{3}{4} \times \dfrac{4}{4} = \dfrac{3}{7} = \underline{0.42}$

Not
Spam

$P(\text{Not spam} \mid X_1 = \text{Yes}, X_2 = \text{Yes}) = P(\text{Not spam}) \times P(X_1 = \text{Yes} \mid NS)$
$$\times P(X_2 = \text{Yes} \mid NS)$$

$$= \frac{3}{7} \times \frac{1}{3} \times \frac{1}{3} = \frac{1}{21}$$

$$= 0.04$$

Compare

$$0.42 > 0.04$$

(spam)    (Not spam)

User new input "spam"
email is a "spam"

# Types

## Gaussian NB

I/P feature are continuous and follow a normal distribution

eg. Iris

## Multinomial NB

I/P features are discrete.

eg. Sentiment Analysis

## Bernoulli NB

I/P features are binary.

eg. Spam detection

| When to Use | When Not to Use |
|---|---|
| * Multi-class classification | * Multicollinearity in data |
| * Simple, fast & easy | * If your data has many features also most of features are irrelevant |
| * Best for small & medium size datasets | |
| * Text Classification | * Imbalanced dataset |
| ⇒ Spam detection | |
| Sentiment Analysis | |