


Naive - Bayes → 70% probability of Rain

→ Proof no hand for similar condition of a day
in past it has rained 70 days out of 100

→ Naive bayes uses the Training data to calculate an observed
probability of each outcome based on evidence given by features

Ex Text classification, Diagnose symptoms.

Information from many features may be considered to increase
the overall performance of model

		↓ CC	↓ Apply	↓ Cuth back	↓ Dredate
Spam ← M ₁ →	Spam	Spam	Spam	Spam	Spam
 ← M ₂ →	Spam	X	Spam	X	X
Spam ← M ₃ →	Spam	Spam	Spam	Spam	Spam

Includes → Weak variables → Impact ↑ if no of Variables ↑

Estimated likelihood of an event is based on the evidence in
Hand across Multiple trials of the event occurrence

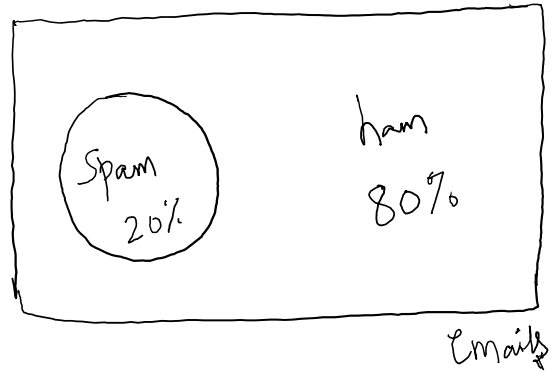
<u>Event</u>	<u>Trial</u>	→ Probability = $\frac{\text{Success trials}}{\text{Total trials}}$
Head	Coin - toss	→ $P(A) = \{0 \text{ to } 1\}$
Spam	Incoming email	→ <u>Mutually exclusive</u> (E) <u>Collectively Exhaustive</u> .
Win lottery	Ticket purchase	

→ Cannot occur at a same time (E) they are the only possible

$$P(\text{spam}) = 0.2$$

$$P(\text{ham}) = 0.8$$

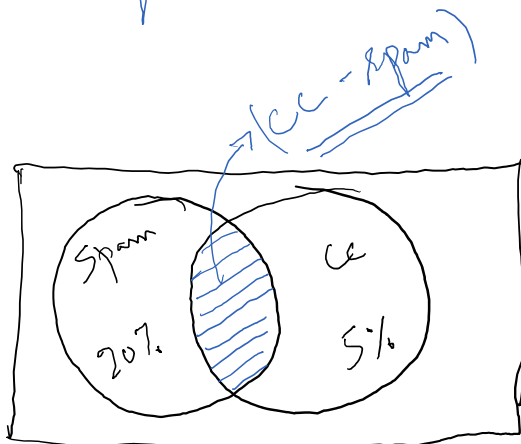
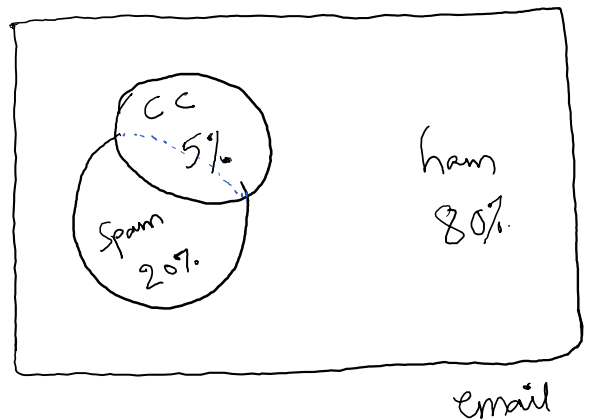
$$P(\text{spam})' = 0.8$$



Joint probability

Not all "cc" mails are spam
(ξ)

Not all spam mails will have "cc"



$P(\text{spam} \cap \text{cc}) \Rightarrow$ Independent event
 \rightarrow Dependent event

Independent event

$$P(A \cap B) = P(A) * P(B)$$

$$= 0.2 * 0.05 = \underline{\underline{0.01}}$$

Dependent event \rightarrow Basis of Predictive modeling

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$P(\text{spam} | \text{cc}) = \frac{P(\text{cc} \cap \text{spam})}{P(\text{cc})} = \frac{P(\text{cc} | \text{spam}) \cdot P(\text{spam})}{P(\text{cc})}$$

Freq table to find components of bayes theorem

	cc		
	Yes	No	
Likelihood			
Spam	4/20	16/20	20
ham	1/20	79/80	80
Total	5/100	95/100	100

$$\begin{aligned}
 P(\text{spam} | cc) &= \frac{P(cc | \text{spam}) * P(\text{spam})}{P(cc)} \\
 &= \frac{4/20 \times 20/100}{5/100} = 0.8
 \end{aligned}$$

For Multiple keywords.

$$P(\text{spam} | w_1, w_2, w_3, w_4) = \frac{P(w_1, w_2, w_3, w_4 | \text{spam}) * P(\text{spam})}{P(w_1, w_2, w_3, w_4)}$$

⇒ Laplace estimation → 20 make it 24

Adding a small no during calculation to Avoid 0

Steps → Building model.

① → Data collection

② → Data Exploration → tm package → (and, but, or) → Removes

→ Split the text into words → "Tokenization"

↪ DTM → Document Text Matrix

→ Visualise → Word cloud (frequency of word occurrence)

→ Indicator creation for frequent words

S₃ → Model → `m ← naivebayes(train=..., class=..., laplace=0)`

$S_4 \rightarrow$ Evaluate \rightarrow Confusion matrix \rightarrow other ways

$S_5 \rightarrow$ Improvement \rightarrow laplace = 1 / change & Run

$S_6 \rightarrow$ Summary \rightarrow $\begin{cases} \rightarrow \% \text{ Accuracy} = ? \\ \rightarrow \text{laplace estimate} = ? \end{cases}$
