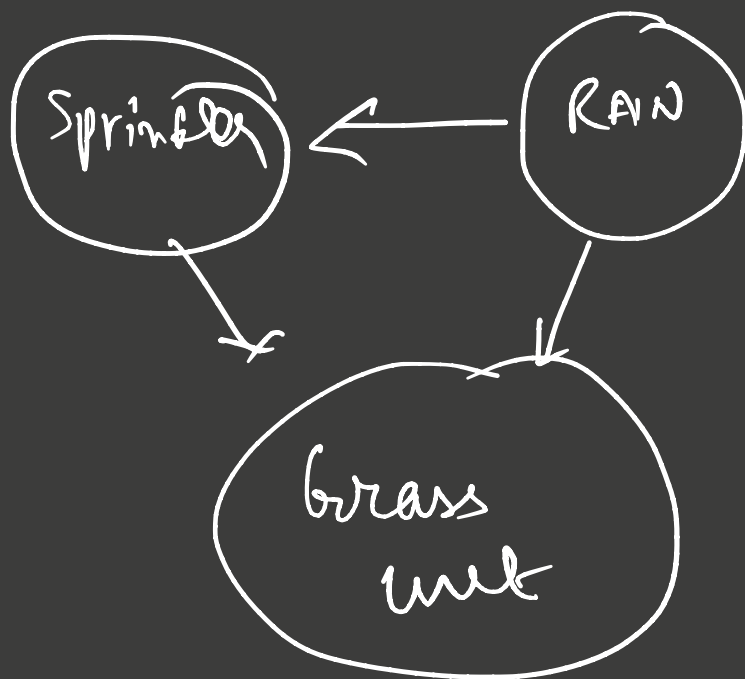


Bayesian Networks

Nodes: Random Variables

Edges: Direct Impact



Classic example

- ① grass could be wet due to
 - RAIN
 - SPRINKLER
- ② If it RAINS, SPRINKLER MAY NOT BE USED

← Random Variables

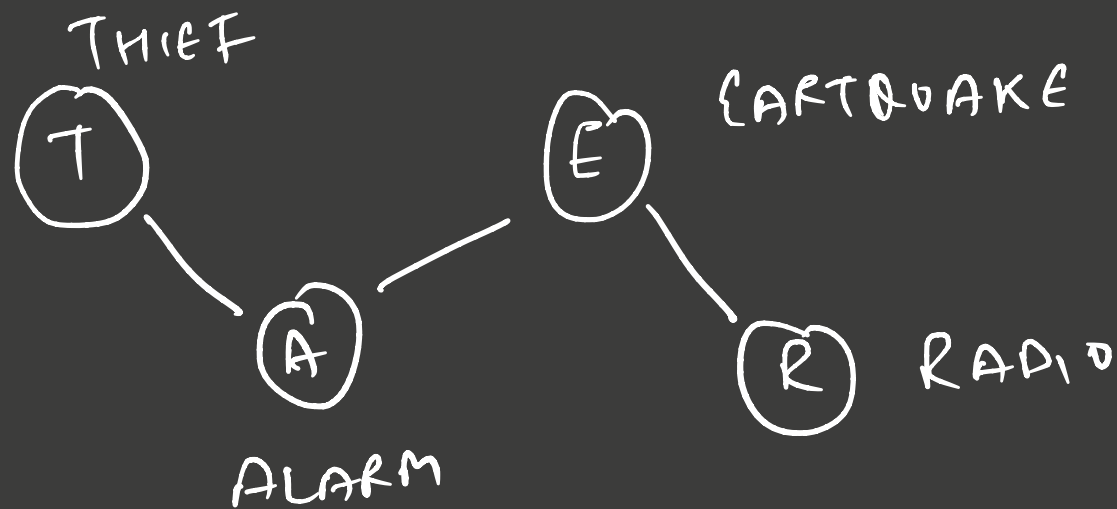
$$P(x_1 \dots x_n) = \prod_{k=1}^n P(x_k | \text{Parents}(x_k))$$

JOINT PROBABILITY

$$c. P(S, R, G) = P(G | S, R) P(S | R) P(R)$$

a) what is $P(y = \text{wet}, \text{Rain} = \text{True}, \text{Sprinkler} = \text{True})$

$$= P(G = w | S = T, R = T) P(S = T | R = T) P(R = T)$$



known
 $P(T), P(E), P(A|T, E), P(R|E)$

$$\stackrel{\text{①}}{=} P(A|T) = \frac{P(A, T)}{P(A)} = \frac{P(A, T, E) + P(A, T, \bar{E})}{P(A, T, E) + P(A, T, \bar{E}) + P(\bar{A}, T, E) + P(\bar{A}, T, \bar{E})}$$

Medical Diagnosis

- (1) you tested +ve for a disease
- (2) Test is 99.1% accurate $\Rightarrow P(\text{Test} = +ve \mid \text{Disease} = \text{True}) = .99$
 $= P(\text{Test} = -ve \mid \text{Disease} = \text{False})$
- (3) Rare disease (1 in 10,000 people)

Q) what is probability you have disease?

From (2) $P(T \mid D) = .99$; $P(\bar{T} \mid \bar{D}) = .99 \Rightarrow P(T \mid \bar{D}) = .01$

From (3) $P(D) = 10^{-4}$; $P(\bar{D}) = 1 - 10^{-4}$

Q is: $P(D \mid T) = ?$ $P(D \mid T) = \frac{P(T \mid D) P(D)}{P(T)}$

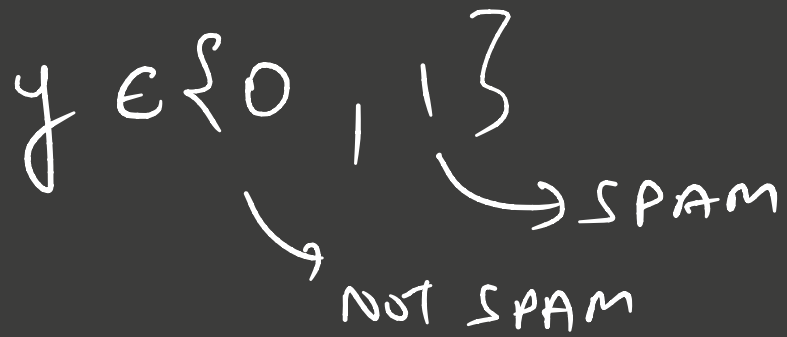
$$P(D|T) = \frac{P(T|D) P(D)}{P(T)}$$

$$= \frac{P(T|D) P(D)}{P(T|D) P(D) + P(T|\bar{D}) P(\bar{D})}$$

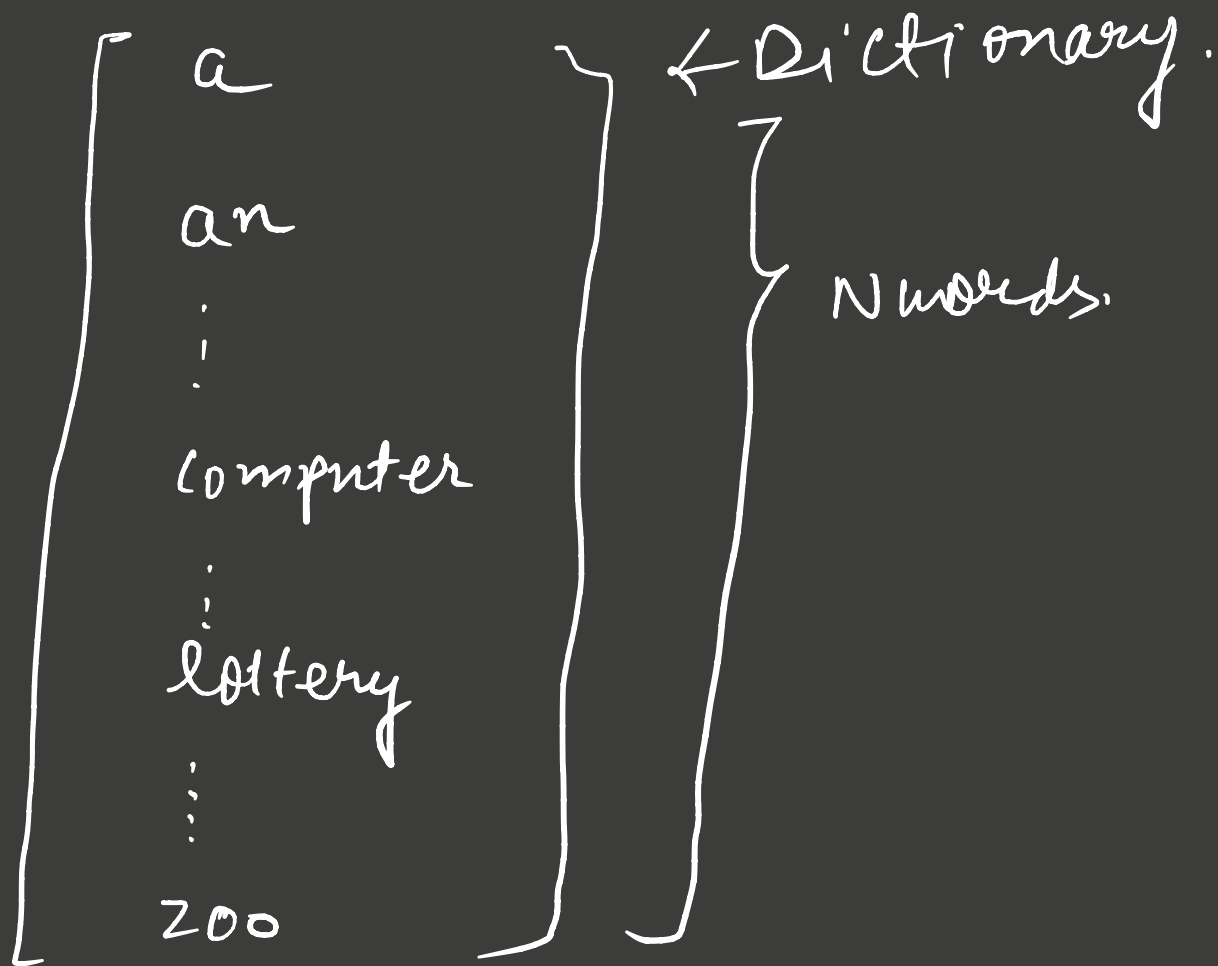
$$= \frac{(.99)(10^{-4})}{(.99)(10^{-4}) + (.01)(1 - 10^{-4})}$$

$$= \frac{(.99)(.0001)}{(.99)(.0001) + (.01)(.9999)} \ll .99$$

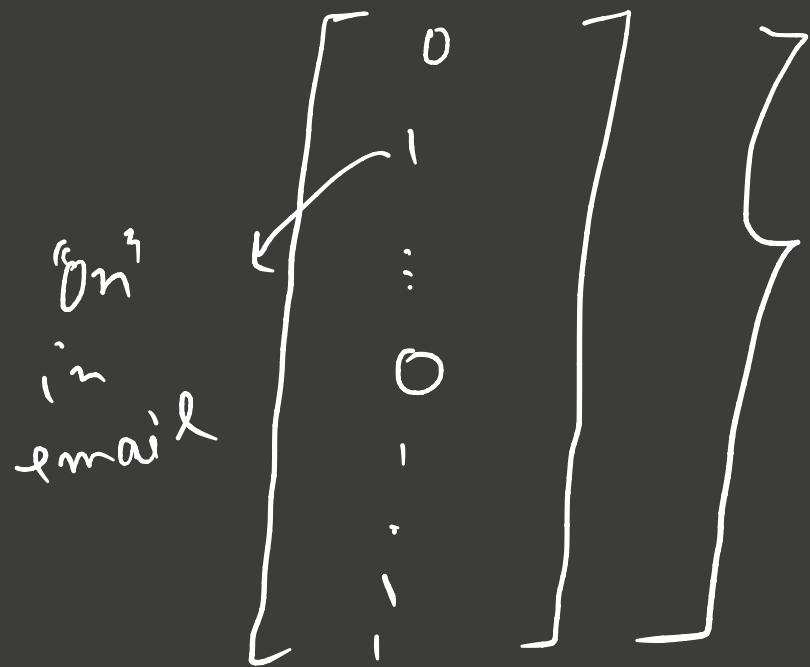
SPAM EMAIL CLASSIFICATION (USING NAIVE BAYES)



WORDS FROM ALL EMAILS



New email.
Construct vector x

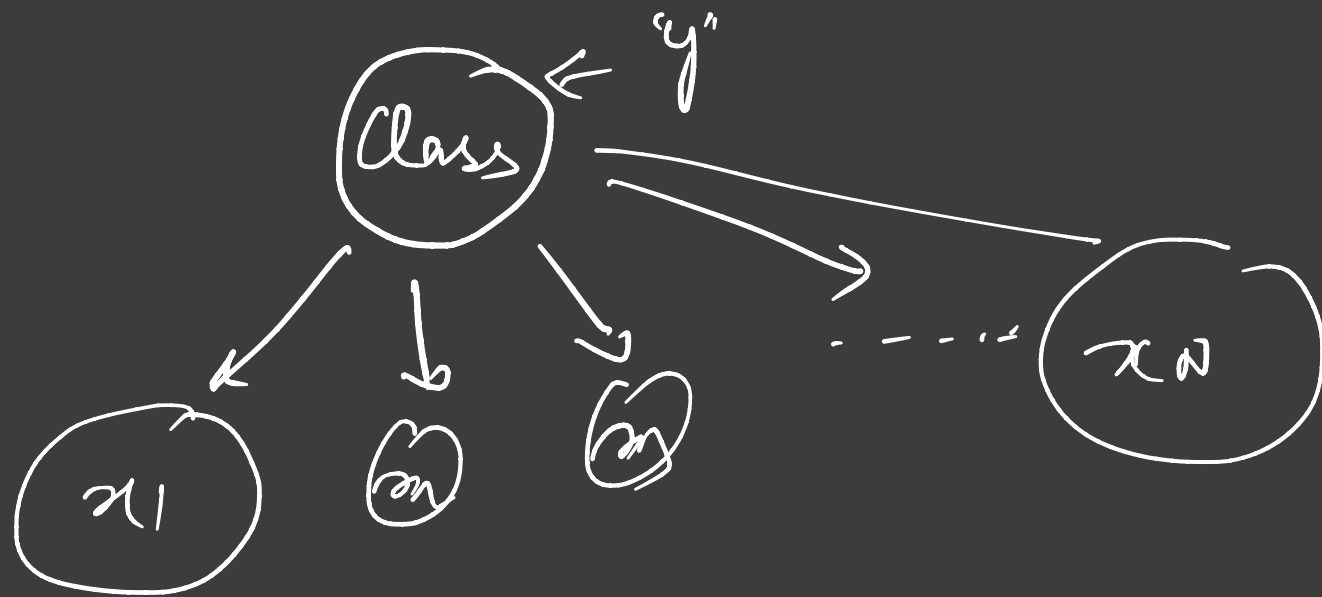


Naive Bayes

* Classification method.

* Scalable.

* Generative.



$$P(x_1, x_2, \dots, x_N | y) = P(x_1 | y) P(x_2 | x_1, y) \dots P(x_N | x_1, \dots, x_{N-1}, y)$$

Naive: Naive assumption x_i and x_{i+1} are independent given y

$$\text{i.e. } P(x_2 | x_1, y) = P(x_2 | y)$$

o
o

$$P(x_1, x_2, \dots, x_n | y) = P(x_1 | y) P(x_2 | y) \dots P(x_n | y)$$

0) what do we want to predict?

$$P(y | x_1, x_2, \dots, x_n)$$

$$= \frac{P(x_1, x_2, \dots, x_n | y) P(y)}{P(x_1, x_2, \dots, x_n)}$$

Parameters

⇒ Probability of x_i being in a spam email

$$P(x_i = 1 | y = 1) = \frac{\text{Count}(x_i = 1 \text{ and } y = 1)}{\text{Count}(y = 1)}$$

⇓

$$P(x_i = 0 | y = 1) = \frac{\text{Count}(x_i = 0 \text{ and } y = 1)}{\text{Count}(y = 1)}$$

⇒ Probability of an email being spam / non-spam

$$P(y = 1) = \frac{\text{Count}(y = 1)}{\text{Count}(y = 1) + \text{Count}(y = 0)}$$

Example

Dictionary is: $[w_1 w_2 w_3]$

TRAIN SET

	w_1	w_2	w_3	y
1	0	0	0	1
2	0	0	0	0
3	0	0	0	1
4	1	0	0	0
5	1	0	1	1
6	1	1	1	0
7	1	1	1	1
8	1	1	0	0
9	0	1	1	0
10	0	0	1	1

	w_1	w_2	w_3	y
1	0	0	0	1
2	0	0	0	0
3	0	0	0	1
4	1	0	0	0
5	1	0	1	1
6	1	1	1	0
7	1	1	1	1
8	1	1	0	0
9	0	1	1	0
10	0	0	1	1

$$y = 0$$

$$P(w_1 = 0 | y = 0) = 3/5 = 0.6$$

$$P(w_2 = 0 | y = 0) = 2/5 = 0.4$$

$$P(w_3 = 0 | y = 0) = 3/5 = 0.6$$

$$P(y = 0) = 0.5$$

$$y = 1$$

$$P(w_1 = 1 | y = 1) = 2/5 = 0.4$$

$$P(w_2 = 1 | y = 1) = 1/5 = 0.2$$

$$P(w_3 = 1 | y = 1) = 3/5 = 0.6$$

$$P(y = 1) = 0.5$$

Test email $\{0, 0, 1\}$

$$P(y=1 | \omega_1=0, \omega_2=0, \omega_3=1)$$

$$= \frac{P(\omega_1=0 | y=1) P(\omega_2=0 | y=1) P(\omega_3=1 | y=1) * P(y=1)}{P(\omega_1=0, \omega_2=0, \omega_3=0)}$$

$$= \frac{(1 - P(\omega_1=1 | y=1)) (1 - P(\omega_2=1 | y=1)) (P(\omega_3=1 | y=1)) * P(y=1)}{2}$$

$$= \frac{(0.6) * (0.8) * (0.6) * 0.5}{2}$$

Similarly,

$$P(y=0 | \omega_1=0, \omega_2=0, \omega_3=1) = \frac{(0.6) * (0.4) * (0.6) * 0.5}{2}$$

$$\frac{P(y=1 | w_1=0, w_2=0, w_3=1)}{P(y=0 | w_1=0, w_2=0, w_3=1)} = \frac{6 \times 8 \times 6 \times 5}{6 \times 4 \times 6 \times 5} = 2$$

$$\therefore P(y=1, \dots) > P(y=0, \dots)$$

\therefore we estimate this to be a "spam"
message

Gaussian Naive Bayes

Classes C_1, \dots, C_k

Continuous attribute x

For class C_k , $\mu_k = \text{Mean}(x | y(x) = C_k)$
 $\sigma_k^2 = \text{Variance}(x | y(x) = C_k)$

Now for $x = \text{some observation 'v'}$

$$P(x=v | C_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(v-\mu_k)^2}{2\sigma_k^2}}$$

Example (from wikipedia)

Height (Feet)	Weight (lbs)	Foot size (inches)	Gender
6	180	12	M
5.92	190	11	M
5.58	170	12	M
5.92	165	10	M
5	100	6	F
5.5	150	8	F
5.42	130	7	F
5.75	150	9	F

	Male	Female
Mean (height)	5.855	5.4175
Variance (height)	3.5×10^{-2}	9.7×10^{-2}
Mean (weight)	176.25	132.5
Variance (weight)	1.22×10^2	5.58×10^2
Mean (foot)	11.25	7.5
Variance (foot)	9.7×10^{-1}	1.67

0) Classify person height = 6ft, weight = 130 lbs, feet = 8 inches.

$$P(\text{male} \mid h=6\text{ft}, \text{weight}=130\text{ lbs}, \text{feet}=8\text{ inches})$$

$$= P(\text{male}) \times P(h=6\text{ft} \mid \text{male}) \times P(\text{weight}=130 \text{ lbs} \mid \text{male}) \times$$

$$P(\text{feet}=8 \mid \text{male})$$

$$= 0.5 \times \frac{1}{\sqrt{2\pi \text{Var}(\text{height} \mid \text{male})}} e^{-\frac{(-6 - 5.833)^2}{2 \text{Var}(\text{height} \mid \text{male})}} \dots$$

$$= \frac{6.2 \times 10^{-9}}{2}$$

$$P(\text{female} | \dots) = \frac{5.4 \times 10^{-7}}{2}$$

\therefore Classified as female.

generating Data | Sampling from Naive Bayes Model,