

### Naive Bayes classifier

$$p(A/B) = p(B/A) p(A) / p(B)$$

$p(\text{disease}/+\text{ve})$  where  $p(\text{disease}) = 1\%$  ,  $p(\text{test+Ve/disease}) = 99\%$   
 $p(\text{disease}/+\text{ve}) = 0.5$

$p(\text{disease}/+\text{ve})$  where  $p(\text{disease}) = 10\%$  ,  $p(\text{test+Ve/disease}) = 99\%$   
 $p(\text{disease}/+\text{ve}) = ??$

$$\begin{aligned} p(\text{disease}/+\text{ve}) &= p(+\text{ve/disease}) p(\text{disease}) / p(+\text{ve}) \\ &= 0.99 * 0.1 / p(\text{test+ve/disease}) p(\text{disease}) + p(\text{test+ve/} \\ &\text{nodisease}) * p(\text{nodisease}) \\ &= 0.99 * 0.1 / ( (0.99 * 0.1) + (0.01 * 0.9)) \\ &= 0.99 * 0.1 / (0.099 + 0.099) \\ &= 0.099 / 0.108 \\ &= 0.916 \end{aligned}$$

### Problem#1: Email classifier using a naive Bayes classifier

training set:

E1 : "good".	not spam (NS)
E2 : " very good"	not spam
E3: "Bad"	spam (S)
E4 : " very bad"	spam
E5: " very bad, very bad"	spam

$$\begin{aligned} p(\text{NS}) &= n(\text{NS}) / n(\text{NS}) + n(\text{S}) = 2 / 2+3 = 2/5 = 0.4 \\ p(\text{S}) &= n(\text{S}) / n(\text{NS}) + n(\text{S}) = 3/5 = 0.6 \end{aligned}$$

unique words are seen in email = {"good", "very", "bad"}

$$\begin{aligned} p(\text{good}) &= 2/5 = 0.4 \\ p(\text{very}) &= 3/5 = 0.6 \\ p(\text{bad}) &= 3/5 = 0.6 \end{aligned}$$

$$p(\text{NS/very}) = 1/3 = 0.33$$

$$p(\text{S/good}) = 0 = 0.05$$

$$p(\text{NS/bad}) = 0 = 0.05$$

$$p(\text{good/NS}) = p(\text{NS/good}) P(\text{good}) / p(\text{NS}) = 0.4 * 0.4 / 0.4 = 0.4$$

$$p(\text{good/S}) = p(\text{S/good}) p(\text{good}) / p(\text{S}) = 0.05 * 0.4 / 0.6 = 0.033$$

$$p(\text{very/NS}) = p(\text{NS/very}) p(\text{very}) / p(\text{NS}) = 0.33 * 0.6 / 0.4 = 0.495$$

$$p(\text{very/S}) = p(\text{S/very}) p(\text{very}) / p(\text{S}) = 0.66 * 0.6 / 0.6 = 0.66$$

$$p(\text{bad/NS}) = p(\text{NS/bad}) * p(\text{bad}) / p(\text{NS}) = 0.05 * 0.6 / 0.4 = 0.075$$

$$p(\text{bad/S}) = p(\text{S/bad}) * p(\text{bad}) / p(\text{S}) = 1 * 0.6 / 0.6 = 1$$

test data

E6 : good bad very bad

$$p(\text{E6/S}) = P(\text{good/S}) * p(\text{bad/S}) * p(\text{very/S}) = 0.033 * 1 * 0.66 = 0.02178$$

$$p(\text{E6/NS}) = p(\text{good/NS}) * p(\text{bad/NS}) * p(\text{very/NS}) = 0.4 * 0.075 * 0.495 = 0.01485$$

$$p(\text{S/E6}) = P(\text{E6/S}) * P(\text{S}) / p(\text{E6}) = 0.6 * p(\text{E6/S}) / p(\text{E6}) = 0.6 * 0.02178 / p(\text{E6}) = 0.0130/p(\text{E6})$$

$$p(\text{NS/E6}) = p(\text{E6/NS}) p(\text{NS}) / p(\text{E6}) = 0.4 * p(\text{E6/NS}) / p(\text{E6}) = 0.4 * 0.01485 / p(\text{E6}) = 0.0059/p(\text{E6})$$

$$\text{say } p(\text{E6}) = 1\% = 0.01$$

$$p(\text{S/E6}) = 0.0013$$

$$p(\text{NS/E6}) = 0.00059$$

since  $p(\text{S/E6}) > p(\text{NS/E6})$  , E6 will be classified as "S" aka SPAM

Problem#2: Classify fruits using a naive Bayes classifier

Let's say the fruit basket consists of fruits as shown in below table

Fruit	Total count
Bananas (B)	50
Orange (O)	30
Other fruit (OF)	100

Attributes	Long (L)	Not long (NL)	Sweet(S)	Not Sweet(NS)	Yellow(Y)	Not yellow (NY)
Bananas	40	10	35	15	45	5
Orange	0	30	15	15	30	0
Other fruit	50	50	65	35	80	20
Total	90	90	115	65	155	25

$$p(B) = n(B) / n(\text{Total fruits in the basket}) = 50 / 50+30+100 = 50/180 = 0.2778$$

$$p(O) = 30/180 = 1/6 = 0.1667$$

$$p(OF) = 100/ 180 = 10/18 = 5/9 = 0.5556$$

$$p(L) = n(L)/n(L) + n(NL) = 90 / 90+90 = 0.5$$

$$p(S) = n(S) / n(S) + n(NS) = 115 / 115+65 = 115/180$$

$$p(Y) = n(Y) / n(Y) + n(NY) = 155 / 155+25 = 155/ 180$$

$$p(L, S, Y) = p(L) * p(S) * p(Y) = 0.5 * 115/180 * 155/180 =$$

$$p(L, S, Y/ B) = p(L/B) * p(S/B) * p(Y/B) = 40/50 * 35/50 * 45/50 =$$

$$p(B/ L, S, Y) = p(L, S, Y/ B) p(B)/ p(L, S, Y) = 40/50 * 35/50 * 45/50 * 0.2778 / (0.5 * 115/180 * 155/180) =$$

$$p(B/ L, S, NY) =$$

$$p(B/ L, NS, Y) =$$

$$p(B/ L, NS, NY)$$

$$p(B/ NL, S, Y) =$$

$$p(B/ NL, S, NY) =$$

$$p(B/ NL, NS, Y) =$$

$$p(B/ NL, NS, NY) =$$

$p(O/L, S, Y) =$   
 $p(O/L, S, NY) =$   
 $p(O/L, NS, Y) =$   
 $p(O/L, NS, NY)$   
 $p(O/NL, S, Y) =$   
 $p(O/NL, S, NY) =$   
 $p(O/NL, NS, Y) =$   
 $p(O/NL, NS, NY) =$

$p(OF/L, S, Y) =$   
 $p(OF/L, S, NY) =$   
 $p(OF/L, NS, Y) =$   
 $p(OF/L, NS, NY)$   
 $p(OF/NL, S, Y) =$   
 $p(OF/NL, S, NY) =$   
 $p(OF/NL, NS, Y) =$   
 $p(OF/NL, NS, NY) =$

Given it is long, sweet & yellow, which fruit it is likely to be ?