

## Naive Bayes

→ It's a probability based classification algorithm.

### Conditional probability

$$P(A|B) = P(A=a | B=b)$$

Probability that  $A=a$  given  $B=b$

$A, B \rightarrow$  Random variables

Ex: let,  $D_1$ : the outcome of die 1

$D_2$ : the outcome of die 2

Question: What is the probability that  $D_1=2$  given that  $D_1 + D_2 \leq 5$

Soln.:  $P(D_1=2 | D_1 + D_2 \leq 5)$

① possible values for  $D_1 + D_2 \leq 5$

$$= \{ (1,1), (1,2), (1,3), (1,4), (2,1), (2,2), (2,3), (3,1), (3,2), (4,1) \}$$

$= 10$

outcomes when  $D_1=2$

$$\Rightarrow P(D_1=2 | D_1 + D_2 \leq 5) = \frac{3/36}{10/36} = \frac{3}{10}$$

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

} Valid only if  $P(B) \neq 0$

Note:  $P(A \cap B) \rightarrow$  probability of both events occurring

$P(A|B) \rightarrow$  Individual probability of  $A$  occurring given that  $B$  has occurred

Q2: Suppose an individual applying to college determines that he has 80% chance of being accepted & he knows that dormitory housing will only be provided to 60% of all the accepted students. Then

$P(\text{Accepted and Dormitory housing})$   
 $= P(\text{Dormitory housing} | \text{Accepted}) * P(\text{Accepted})$   
 $= 0.6 * 0.8 = 0.48$

Independent v/s Mutually Exclusive events

→ A & B are independent }  $P(A|B) = P(A)$   
 $P(B|A) = P(B)$

Example 1:

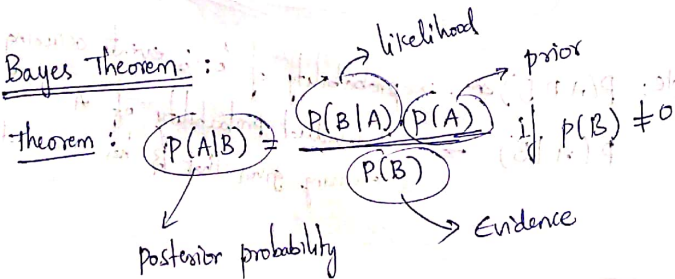
A: Getting value of 6 in die 1 throw ( $D_1=6$ )  
 B: " " " 3 in die 2 throw ( $D_2=3$ )

⇒  $P(D_1=6 | D_2=3) = P(D_1=6) = 1/6$

→  $P(A|B) = P(B|A) = 0$  } then A & B are mutually exclusive events

$\frac{P(A \cap B)}{P(B)} = \frac{P(B \cap A)}{P(A)} = 0$  }  $A \cap B = B \cap A = 0$  for mutually exclusive events

Bayes Theorem:



proof:

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

$A \cap B = B \cap A \rightarrow$  from Set theory

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

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

⇒  $P(B \cap A) = P(B|A) * P(A)$

∴  $P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$

Naive Bayes Algorithm

Naive Bayes is Bayes theorem  
 unsophisticated or simplistic

→  $P(C_k | x_1, x_2, \dots, x_n)$   
 class label | vector of features (independent) of a data point

$P(C_k | x) = \frac{P(C_k) \cdot P(x | C_k)}{P(x)}$

$= \frac{P(C_k) \cdot P(x)}{P(x)}$

posterior =  $\frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$



$$\rightarrow P(C_k | x) = \frac{P(C_k) \cdot P(x | C_k)}{P(x)}$$

In practice, there is interest only in the numerator of the above fraction, because the denominator  $P(x)$  is  $x$  independent of  $C_k$ .

$\therefore$  we consider only numerator

$$\Rightarrow P(C_k | x) = P(C_k) \cdot P(x | C_k) = P(x \cap C_k) = P(C_k, x)$$

Joint probability model

$\rightarrow$  ~~for  $C_1$~~   $P(C_1 | x)$  is high for class " $1$ "  
then we consider " $x$ " to belong to class " $1$ "

$$P(C_k, x) = \text{Joint probability model}$$

$$\rightarrow P(C_k, x_1, x_2, \dots, x_n) = P(x_1, x_2, \dots, x_n, C_k)$$

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

$$= P(x_1 | x_2, \dots, x_n, C_k) \cdot P(x_2 | x_3, x_4, \dots, x_n, C_k) \cdot P(x_3, \dots, x_n, C_k)$$

$$\vdots$$

$$= P(x_1 | x_2, \dots, x_n, C_k) \cdot P(x_2 | x_3, x_4, \dots, x_n, C_k) \cdot \dots \cdot P(x_{n-1} | x_n, C_k) \cdot P(x_n | C_k) \cdot P(C_k)$$

$\rightarrow$  In naive Bayes, we assume that each feature  $x_i$  is conditionally independent of every other feature

Simple independence:  $P(A|B) = P(A)$

Conditional independence:  $P(A|B, C) = P(A|C)$

$$\Rightarrow P(x_1 | x_{1+1}, \dots, x_n, C_k) = P(x_1 | C_k) \cdot \underbrace{P(x_{1+1}, \dots, x_n | C_k)}_{\text{independent of } x_1 \text{ given } C_k}$$

Note: "Naive" assumption: Features are conditionally independent of each other.

$\rightarrow$  Thus the joint model can be expressed as

$$P(C_k | x_1, \dots, x_n) \propto P(C_k, x_1, \dots, x_n) \propto P(C_k) \cdot P(x_1 | C_k) \cdot P(x_2 | C_k) \cdot \dots \cdot P(x_n | C_k) \propto P(C_k) \cdot \prod_{i=1}^n P(x_i | C_k)$$

$$\Rightarrow P(C_k | x_1, \dots, x_n) = \frac{1}{Z} P(C_k) \cdot \prod_{i=1}^n P(x_i | C_k)$$

where the evidence,  $Z = P(x)$

$$= \sum_K P(C_k) P(x | C_k)$$

is a scaling factor dependent only on  $x_1, \dots, x_n$

$$y = \underset{k \in \{1, 2, \dots, K\}}{\text{argmax}} P(C_k) \cdot \prod_{i=1}^n P(x_i | C_k)$$

# Toy Example: Train and Test Stages

Source: [skitterline.com/naive-bayes](http://skitterline.com/naive-bayes)

Date	Predictors				Response
	outlook	temperature	humidity	wind	play = Yes/No
Day 1	Sunny	hot	high	weak	Yes
Day 2	overcast	mild	Normal	strong	No
Day 3	Rain	cool	Normal	weak	Yes

the learning (or training) phase

In the learning phase, we compute the table of likelihoods (probabilities) from the training data

→ In Naive Bayes,

$$P(C|f_1, f_2, f_3, f_4) \propto P(f_1|C) \cdot P(f_2|C) \cdot P(f_3|C) \cdot P(f_4|C) \cdot P(C)$$

$$P(\text{outlook} = \text{sunny} | C = \text{yes})$$

likelihoods

outlook (f <sub>1</sub> )	Frequency		Probabilities	
	play = Yes	play = No	P(f <sub>1</sub>  C=Yes)	P(f <sub>1</sub>  C=No)
Sunny	2	3	2/9	3/9
Overcast	4	0	4/9	0/5
Rain	3	2	3/9	2/5
	9	5		

→ likelihood for f<sub>1</sub> in outlook

Similarly, calculate the likelihood table for other features: temp, humidity and wind.

→ Also calculate P(C)

$$P(C = \text{Yes}) = 9/14$$

$$P(C = \text{No}) = 5/14$$

Time complexity: O(N \* D) with optimization it is O(N \* D)

## The Classification phase

Let's say, we get a new instance of the weather condition, ~~for~~

x' = (outlook = sunny, Temp = cool, Humidity = high, wind = strong)

Question: Are we going to play tennis under the conditions specified by x'?

Soln:

$$P(\text{play} = \text{yes} | x') \propto P(\text{Sunny} | \text{play} = \text{yes}) \times P(\text{cool} | \text{play} = \text{yes}) \times P(\text{high} | \text{play} = \text{yes}) \times P(\text{strong} | \text{play} = \text{yes}) \times P(\text{yes})$$

$$\propto \frac{2}{9} \times \frac{1}{9} \times \frac{3}{9} \times \frac{1}{9} \times \frac{9}{14}$$

$$\propto \frac{1}{189} \propto 0.0053$$

$$P(\text{play} = \text{no} | x') \propto P(\text{outlook} = \text{sunny} | \text{play} = \text{no}) \times P(\text{temp} = \text{cool} | \text{play} = \text{no}) \times P(\text{hum} = \text{high} | \text{play} = \text{no}) \times P(\text{wind} = \text{strong} | \text{play} = \text{no}) \times P(\text{play} = \text{no})$$

$$\propto \frac{3}{9} \times \frac{1}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{5}{14} \propto 0.0205$$

Since P(play=no|x') > P(play=yes|x'), we classify the new instance x' to be the conditions where game cannot be played, i.e. for x', we give y' = 0

NOTE: The above implementation is for categorical features



# Naive Bayes on Text data

→ Naive Bayes is very popular for Text classification

→ Example:

	class
test 1 →	0
test 2 →	1
test 3 →	0
...	...

Task  
 $P(y_q = 1 | x_q)$ ?

$P(y_q = 0 | x_q)$ ?

text → preprocessing →  $w_1, w_2, \dots, w_n$  → Binary Bow

$P(y=1 | \text{test}) = P(y=1 | \underbrace{w_1, w_2, \dots, w_n}_{\text{features}})$

$\propto P(y=1) \times P(w_1 | y=1) \times P(w_2 | y=1) \dots P(w_n | y=1)$

$\propto P(y=1) \prod_{i=1}^d P(w_i | y=1)$

Similarly

$P(y=0 | \text{test}) \propto P(y=0) \prod_{i=1}^d P(w_i | y=0)$

→  $P(y=1) = \frac{\# \text{training pts with } y=1}{\text{total } \# \text{ training pts}}$

$P(y=0) = \frac{\# \text{training pts with } y=0}{\text{total } \# \text{ training pts}}$

→ To compute the likelihood probabilities, divide the training set into two parts ( $y=1$  &  $y=0$ )

→  $P(w_i | y=1) = \frac{\# \text{data pts with } w_i \& y=1}{\# \text{data pts with } y=1}$

NOTE: For text-classification problems, Naive Bayes is a very good baseline.

## Laplace / Additive Smoothing

Training

$\left. \begin{aligned} &P(y=1) ; P(y=0) \leftarrow \text{class priors} \\ &P(w_1 | y=1) ; P(w_1 | y=0) \\ &P(w_2 | y=1) ; P(w_2 | y=0) \\ &P(w_n | y=1) ; P(w_n | y=0) \end{aligned} \right\} \text{likelihoods}$

Test: ~~test~~  $\text{test}_q = (w_1, w_2, w_3, \dots, w_n)$   
 Say,  $w^1$  is not present in list of words of  $w_1, w_2, \dots, w_n$   
 then what to do?

$P(y=1 | \text{test}_q) = P(y=1 | w_1, w_2, w_3, w^1)$

$= P(y=1) \times P(w_1 | y=1) \times P(w_2 | y=1) \times P(w_3 | y=1) \times P(w^1 | y=1)$

How do we get this?

→  $P(w^1 | y=1) = \frac{P(w^1, y=1)}{P(y=1)}$

$= \frac{\# \text{pts such that } w^1 \text{ occurs in } y=1}{\# \text{pts where } y=1}$

$= \frac{0}{n_1} = 0$

Note If we substitute above  $P(w' | y=1) = 0$ , then we get  $P(y=1 | \text{test})$  also zero. That's a problem! There should be some way of handling it.

The workaround for this problem is Laplace/Additive Smoothing.

$$\rightarrow P(w' | y=1) = \frac{0 + \alpha}{n_1 + \alpha K}$$

$\alpha = 1$  typically in this case  $K=2$   
 $K = \# \text{ distinct values } w' \text{ can take}$

$$\rightarrow \text{let } n_1 = 0$$

$$\Rightarrow P(w' | y=1) = \frac{0 + \alpha}{100 + 2\alpha}$$

Case 1:  
 $\alpha = 1 = \frac{1}{102} \neq 0 \Rightarrow P(w' | y=1) \neq 0$

Why  $\frac{0 + \alpha}{n_1 + \alpha K}$ ?

Case 2:  $\alpha = 10000$

$$\Rightarrow P(w' | y=1) = \frac{0 + 10000}{100 + 20000} = \frac{10000}{20100} \approx \frac{1}{2}$$

when  $\alpha$  is very large (like in case 2)

we get

$$P(w' | y=1) = P(w' | y=0) = \frac{1}{2}$$

$\rightarrow$  In Laplace Smoothing, we add  $\alpha$  &  $\alpha K$  for all the words.

$$\therefore P(w_i | y=1) = \frac{\# \text{ data pts with } w_i \& y=1 + \alpha}{\# \text{ data pts } y=1 + \alpha K}$$

Note:  $\alpha \uparrow \Rightarrow$  moving likelihood probabilities to uniform distribution

The prob will move close to  $1/2$

Since we are adding terms in numerator & denom, it is called additive Smoothing.

$\Rightarrow$  often times,  $\alpha = 1$ , this is called add-one Smoothing.

Log-probabilities for numerical stability

We have

$$P(y=1 | w_1, w_2, w_3, \dots, w_d)$$

$$= P(y=1) \times P(w_1 | y=1) \times P(w_2 | y=1) \dots \times P(w_d | y=1)$$

Here all probabilities lie b/w 0 & 1.

If we let  $d = 100$ , then multiplying 100 small values ( $0 < p < 1$ ), will give very small value. This causes

numerical stability issues. (ex: numerical underflow)

$\rightarrow$  One soln for this is to use "log" probabilities



## Log probabilities

$$P(y=1 | w_1, w_2, \dots, w_d) = P(y=1) \times \prod_{i=1}^d P(w_i | y=1)$$

$$P(y=0 | w_1, w_2, \dots, w_d) = P(y=0) \times \prod_{i=1}^d P(w_i | y=0)$$

→ 'log' is monotonic func.  $\Rightarrow$   $\uparrow$ ;  $\log \uparrow$

$$\log(P(y=1 | w_1, w_2, \dots, w_d))$$

$$= \log(P(y=1)) + \sum_{i=1}^d \log(P(w_i | y=1))$$

## Bias and Variance tradeoff

High bias  $\rightarrow$  Underfitting  
High variance  $\rightarrow$  Overfitting

→ For Naive Bayes,  $\alpha$  in Laplace Smoothing determines underfit (or) overfit  
high bias  $\rightarrow$  high variance

→ Case 1:  $\alpha=0$

$$P(w_i | y=1) = \frac{\text{train data points where } w_i \text{ occurs } y=1}{\text{train pts with } y=1}$$

Even for words which are very rare, we're giving probabilities  
This results in overfit

when  $\alpha=0 \Rightarrow$  Small change in train results in large change in the model  $\Rightarrow$  high variance  $\Rightarrow$  overfitting

→ Case 2: when  $\alpha$  is large,

$$P(w_i | y=1) \approx P(w_i | y=0) \approx \frac{1}{2}$$

$$\therefore P(y=1 | w_1, w_2, \dots, w_d) = P(y=1) \times \prod_{i=1}^d P(w_i | y=1)$$

$$P(y=0 | w_1, w_2, \dots, w_d) = P(y=0) \times \prod_{i=1}^d P(w_i | y=0)$$

$\Rightarrow$   $n_1$  (true pts)  $\approx n_2$  (false pts)

$$\text{Since } n_1 > n_2 \Rightarrow P(y=1) > P(y=0)$$

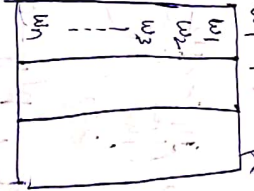
$\therefore P(y=1 | w_1, w_2, \dots, w_d) > P(y=0 | w_1, w_2, \dots, w_d)$   
For most of the points, it will classify as true  
This is underfitting

Question - how to find right " $\alpha$ "?  
Ans: Using Simple cross validation or 10-fold cross validation

$\alpha$  in NB is hyperparameter

## Feature importance & Interpretability

We have,  $w_1, P(w_1 | y=1), P(w_1 | y=0)$



① Sorts  $w_j$ 's based on  $P(w_j | y=1)$

Say in descending order

② words which have high

$P(w_j | y=1)$  are important

features in determining the point  $\in$  the class

→ Feature importance in NB:  
for +ve class } find  $w_i$  with highest value of  $P(w_i | y=1)$

for -ve class } find  $w_i$  with lowest value of  $P(w_i | y=0)$

Note: In NB, important features are obtained directly from the model.

Interpretability  
 $x_q \rightarrow y_q = 1$   
I am concluding  $y_q = 1$  because  $x_q$  contains words  $w_3, w_6, w_{10}$  which have high value of  $P(w_i | y=1)$

### Imbalanced data

$$P(y=1 | w_1, w_2, \dots, w_d) = P(y=1) \cdot \prod_{i=1}^d P(w_i | y=1)$$

$$P(y=0 | w_1, w_2, \dots, w_d) = P(y=0) \cdot \prod_{i=1}^d P(w_i | y=0)$$

Since imbalanced data }  $\Rightarrow \frac{n_1}{n} \gg \frac{n_2}{n}$   
(pos) (neg)  
90% 10%

Then  $P(y=1) = \frac{n_1}{n} \approx 0.9$ ,  $P(y=0) = 0.1$

when we have imbalanced data then  
① Due to class priors  $\rightarrow$  majority/dominant class has an advantage.

Soln: ① Upsampling (or) Downsampling  
② Drop  $P(y=1)$  &  $P(y=0)$   
 $\rightarrow$  make  $P(y=1) = P(y=0) = \frac{1}{2}$  or  $\frac{1}{4}$

→ For an imbalanced dataset, while applying Laplace Smoothing, the effect of  $\alpha$  is more on minority class probability.

### Outliers

①  $x_q = w_1 w_2 w_3 w$  set of words in

$w \notin \{w_1, w_2, \dots, w_m\}$  Drain

$$P(w | y=1) = \frac{0}{n_1} \Rightarrow w \text{ is an outlier.}$$

$\therefore$  Use Laplace Smoothing.

Note: Outlier occurs very few times in +ve & -ve classes

$\rightarrow$  Drop outliers or apply Laplace Smoothing

### Missing values

- ① Text-data  $\rightarrow$  No case of missing data
- ② Categorical features  $\rightarrow$  considered as category
- ③ Numerical features  $\rightarrow$  Imputation



## Handling Numerical features (Gaussian NB)

### Multi class classification

#### Similarity or distance metrics

→ NB cannot use similarity or distance metrics

#### Large dimensionality

→ NB can be used with large dimensional data when dimensionality is large, we log probabilities

#### Best & worst cases for NB

① Conditional independence of features (assumption)

- ② Used for text classification
- ③ Extensively used when we have categorical features
- ④ Interpretability
- ⑤ Runtime complexity, Train time complexity, Runtime space complexity is low.