

2nd

Naive Bayes  $\Rightarrow$  classification algorithm  
 $\hookrightarrow$  probability-based

KNN  $\rightarrow$  neighborhood based classification.

Conditional Probability  $= (P(A/B)) = P(A=a/B=b)$   
 $\downarrow \qquad \qquad \downarrow$   
value that A takes    value that B takes.

Always read equations in english.

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

\* Independent Events & Mutually Exclusive events

A, B are said to be independent.

$$P(A/B) = P(A)$$

$$P(B/A) = P(B)$$

for ex  $\Rightarrow$

A: getting value of 6 in die 1  
throw ( $D_1 = 6$ )

B: getting a value of 3 in  
die 2's throw ( $D_2 = 3$ )

A, B are said to be mutually exclusive if.

$$P(A/B) = P(B/A) = 0$$

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

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

$\therefore$  so,  $P(A \cap B)$  should be 0

$$\text{as } A \cap B = B \cap A$$

for ex  $\Rightarrow$  probability of getting 3 in die 1 is 0 if die 2  
is getting 6 in it.



Bayes Theorem  $\rightarrow$  (1700s)

$$P(A/B) = \frac{P(B/A) \cdot P(A)}{P(B)} \quad \text{if } P(B) \neq 0.$$

$P(A/B)$  is labeled **posterior probability**  
 $P(B/A)$  is labeled **likelihood**  
 $P(A)$  is labeled **prior**  
 $P(B)$  is labeled **evidence**

Proof of Bayes Theorem  $\rightarrow$

$$P(A/B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A, B)}{P(B)} \quad \text{also means } \cap$$

In set theory  $A \cap B = B \cap A \leftarrow$

$$P(A/B) = \frac{P(B \cap A)}{P(B)} = \frac{P(B, A)}{P(B)} \quad \text{--- (1)}$$

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

$$P(B \cap A) = P(B/A) \cdot P(A) \quad \text{--- (2)}$$

Put  $P(B \cap A)$  in (1) eq

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



## \* Naive Bayes Algorithm

↓  
unobtrusive  
↓  
Simplistic

$$P(C/f_1, f_2, f_3, f_4) = P(f_1/C) * P(f_2/C) * P(f_3/C) * P(f_4/C) * P(C)$$

Naive Bayes is much better in terms of space complexity at runtime

$O(d \times c) \rightarrow$  naive Bayes

$O(n \times d) \rightarrow$  KNN

memory efficient at run-time

## \* Naive Bayes on text data $\rightarrow$

first technique used in case of spam filters  $\rightarrow$  email  $\rightarrow$  spam  
↓  
not spam

review  $\rightarrow$  +ve  
↓  
 $\rightarrow$  -ve

Text classification  $\rightarrow$  Naive Bayes is a simple used method.

← text1 →	0
← text2 →	1
	0
	1

Task

①  $P(y_v = 1 / \text{text}_q) \rightarrow$  compute these

②  $P(y_q = 0 / \text{text}_q)$

So if ① > ②

then we can say that  $y_q = 1$  for given  $\text{text}_q$ .



preprocessing  $\rightarrow$

text  $\rightarrow$

stopwords

stemming

n-grams

$\rightarrow \{w_1, w_2, w_3, \dots, w_d\}$

$\downarrow$

Binary Bow

text  $\xrightarrow{\text{pre-pro}}$   $\{w_1, w_2, \dots, w_d\}$   
set of words.

$$P(y_1 = 1 / \text{text}) = P(y_1 = 1 / \underbrace{w_1, w_2, w_3, \dots, w_d}_{\text{features}})$$

class  
prior

likelihoods

$$= P(y=1) * P(w_1/y=1) * P(w_2/y=1) * \dots * P(w_d/y=1)$$

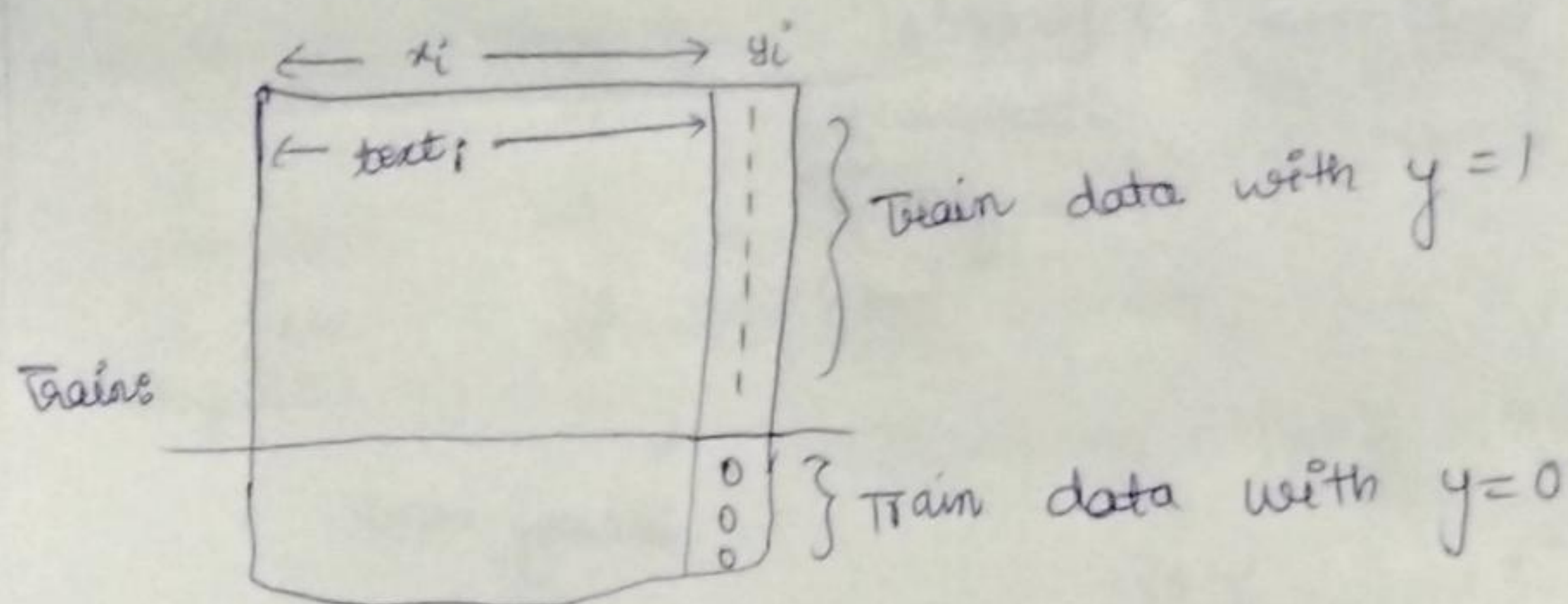
$$P(y=1 / \text{text}) \propto P(y=1) * \prod_{i=1}^d P(w_i / y=1)$$

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

$$P(y=1) = \frac{\text{no. of train points with } y=1}{\text{Total no. of train points.}}$$

$$P(y=0) = \frac{\text{no. of train points with } y=0}{\text{Total no. of train points.}}$$





$$P(w_i | y=1) = \frac{\text{no. of } \overset{\text{Train}}{\text{data points which contain } w_i \&\&}}{\text{no. of } \overset{\text{Train}}{\text{data points with } y=1}}$$

Text - classification problems.

{ spam - detection  
 Reliability of a review } Naive Bayes is a very good Baseline

benchmark

\* Laplace - Additive Smoothing  $\rightarrow$  ~~not Laplace smoothing~~

After training  $\rightarrow$  {  $P(y=1); P(y=0) \leftarrow$  class priors

all this data is already computed. {

$P(w_1/y=1)$     $P(w_1/y=0)$   
 $P(w_2/y=1)$     $P(w_2/y=0)$   
 $\vdots$   
 $P(w_m/y=1)$     $P(w_m/y=0)$

} likelihoods



Test:  $\rightarrow \text{test}_q = (w_1, w_2, w_3, w')$

$\nwarrow$   $w'$  is not present in  $\{w_1, w_2, w_3, \dots, w_m\}$   $\nearrow$  training data.  
very often

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

$$= P(y=1) * P(w_1/y=1) * P(w_2/y=1) * P(w_3/y=1)$$

$$* P(w'/y=1)$$

ignoring or dropping it will mean

$P(w'/y=1) = 1$   
which is not correct.

how do you get <sup>this</sup> probability as  $w'$  is not present in training data.

we have to get values of  $P(w'/y=1)$  and  $P(w'/y=0)$

$$P(w'/y=1) = \frac{P(w', y=1)}{P(y=1)}$$

$$= \frac{\text{no. of train points such that } w' \text{ occurs in } y=1}{\text{no. of train points where } y=1}$$

$$= \frac{0}{n_1} = 0 \rightarrow \text{This is also dangerous as it will make whole probability to be 0.}$$



Laplace smoothing or additive smoothing

$$P(f_i = a / y = 1) = \frac{0 + \alpha}{n_i + \alpha k}$$

$\alpha = 1$  typically (not always)

$k$  = no. of distinct values which  $f_i$  can take

$f_i \Rightarrow$  feature

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

$\uparrow$   
 $k = 2$  because  $w'$  is 0 or 1,  
present or not.

Let  $n_i = 100$

Case 1:  $\alpha = 1 \Rightarrow \frac{1}{102} \neq 0$

$\alpha$  small,  
we are getting  
rid of multiplying  
all the  
probabilities  
with 0.

$$P(w' / y = 1) \neq 0 \quad \Downarrow \text{which implies}$$

$$P(y = 1 | \text{text}_q) \neq 0.$$

So it's not 0 anymore

Case 2:  $\alpha = 10000 \rightarrow$  when  $\alpha$  is large

$$P(w' / y = 1) = \frac{0 + 10k}{100 + 2 \times 10k} = \frac{10k}{20100} \approx \frac{1}{2}$$

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

means equal probability of  $w'$  to be 0 or 1  
because  $w'$  have only two possibilities (0 or 1)



Laplace smoothing  $\Rightarrow$  find this for all words

$$P(w_i | y=1) = \frac{(\text{no. of data points with } w_i \text{ \& } y=1) + \alpha}{(\text{no. of data points with } y=1) + \alpha k}$$

present in  
my training  
data

adding something  
to numerator &  
denominator,  
that's why it is  
called additive  
Smoothing

In this formula, as  $\alpha \uparrow$ ,

$P(w_i | y=1) \rightarrow$  likelihood probability  
is moving to a uniform distribution.

if  $n_i$  is small

num or den is small  $\rightarrow$  less confidence in ratio

often times,  $\alpha = 1 \rightarrow$  add one smoothing

It is called smoothing because you are moving likelihood probability towards uniform distribution

\* log-probabilities for numerical stability  $\Rightarrow$

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

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

$x \uparrow ; \log x \uparrow \rightarrow$  monotonic fun.



## \* Bias and Variance Tradeoff $\Rightarrow$

Case 1  $\Rightarrow \alpha = 0 \rightarrow$  small change in  $D_{\text{train}}$  results in large change in the model.  $\Rightarrow$  high variance = overfitting.

Case 2  $\Rightarrow \alpha = \text{very large} = \alpha = 10000$

$$P(w_i/y=1) = \frac{2 + 10000}{1000 + \underbrace{20000}_{K=2}} \approx \frac{1}{2}$$

$\downarrow$   
 $0 \propto 1$

high bias  $\Leftarrow$  (Underfitting) like in case of KNN ( $K=n$ ) case.

$\downarrow$   
based on  $P(y=1)/P(y=0)$

Q How to find the right  $\alpha$ ?

$\Rightarrow$  using simple cross ~~validation~~ validation or 10-fold cross validation

## \* Feature Importance and Interpretability $\Rightarrow$

words  $w_i$ 's which have high values of  $P(w_i/y=0)$  and  $P(w_i/y=1) \rightarrow$  important/meaningful words (features)

So, unlike KNN, we can directly find important features by using likelihoods in Naive Bayes.

Interpretability  $\Rightarrow$  We can give reason by saying  $\rightarrow$

$$x_q \rightarrow (y_q = 1)$$

as  $x_q$  contains words  $w_3, w_6, w_7 \dots$  etc for which  $P(w_3/y=1), \dots$  are high.



\* Imbalanced data  $\Rightarrow$

### Impact on Priors

In case of Imbalanced data, because of class priors ( $P(y=1)$  and  $P(y=0)$ ) majority/ dominating class has an advantage.

Soln ① upsampling or downsampling.

motive  $\Rightarrow n_1 \approx n_2$  and  $P(y=1) = P(y=0) = \frac{1}{2}$

② drop  $P(y=1)$  &  $P(y=0)$  by putting  
 $P(y=1) = P(y=0) = 1$

③ modified naive Bayes for class-imbalance  
not often used

### Impact on likelihood ratios

minimum (-ve)

$$P(w_i/y=1000) = \frac{2}{100} = 2\%$$

let's assume,  $\alpha = 10$

$$10\% = \frac{12}{120} = \frac{2+10}{100+20}$$

maximum (+ve)

$$\frac{18}{900} = 2\%$$

$$\frac{18+10}{900+20} = \frac{28}{920} = 3.04\%$$

minority class  
label is  
benefited.

which is not good.

Soln ① upsampling or downsampling

② Harker's forex

$\alpha_{+ve}$ ,  $\alpha_{-ve}$

~~two~~  $\alpha$ s



\* Outliers  $\rightarrow$  Both in training & test data  
hack/solution  $\rightarrow$

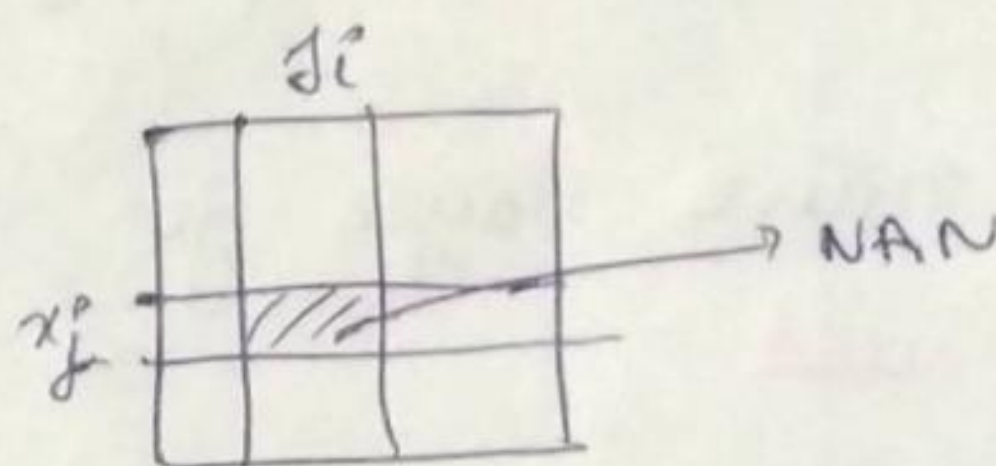
- ① If a word ( $w_j$ ) occurs fewer than  $\tau = 10$  times, then just ignore that word.
- ② Apply Laplace smoothing to that word with a reasonable value of  $\alpha$ .

\* Missing values  $\rightarrow$

① Text-data  $\rightarrow$  no case of missing data, as a word either occurs or not.

② Categorical features  $\rightarrow$

Hack  $\rightarrow$  Consider NAN as a category too.



So,  $f_i \in \{a_1, a_2, a_3, \dots, \text{NAN}\}$

② Numerical features  $\rightarrow$  use of standard feature imputation methods like mean, median etc. based.

Gaussian Naive Bayes.

Naive Bayes with Numerical or realvalued features  $\rightarrow$

assume  $f_j$  in  $D' \rightarrow$  Gaussian dist.

$\downarrow$   
 $D^{+ve}$   
train

$$N(\mu_j^+, \sigma_j^+)$$

assume  $f_j$  in  $D'' \rightarrow$  Gaussian dist

$\downarrow$   
 $D^{-ve}$   
train

$$N(\mu_j^0, \sigma_j^0)$$



This is also called Gaussian Naive Bayes.  $\rightarrow$  likelihood probabilities  
we take  $P(w_i/y=1) \approx P(w_i/y=0)$  values from Gaussian curve, that's why its name is this.

$$P(w_i/y=1)$$

$\swarrow \searrow$   
0 1  $\rightarrow$  Bernoulli Naive Bayes.

another  $\rightarrow$  Multinomial Naive Bayes.  
is

Naive Bayes is making an assumption of conditional independence.

Multiclass Classification  $\rightarrow$

Same as  $P(y_i=0/w_1, w_2, \dots, w_d)$

we can compute as many classes probability

like  $P(y_i=2/w_1, w_2, \dots, w_d)$

and compare them.  $\rightarrow$  Naive Bayes is by default made for multiclass classification.

Similarity matrix or Distance matrix  $\rightarrow$

Naive Bayes cannot use similarity or distance matrix.

as to compute probabilities (likelihood) we need actual features values for that.



#### \* Large Dimensionality →

Naive Bayes <sup>is</sup> easily applied on large dimensionality as in case of text classification.

Just use log-probabilities in place of normal probabilities to maintain numerical stability.

#### \* Best and worst cases for Naive Bayes →

① Naive Bayes perform very well with conditional independence assumption, but conditional independence ↓ NB degrades.

But when some features are dependent → NB works reasonably well.

→ good  
→ great  
→ Phenomenal } → conditional dependent.

② Text Classification → email spam  
↓  
review polarity { high-dim }

NB is baseline for these problems.

③ Categorical features → NB is used  $\Leftarrow$  much.

④ Real-valued features → NB is not often used.

⑤ Interpretable, feature importance →

run-time complexity → low  
train time complexity → low  
run-time space → low.

medical domain



NB → Simple to implement  
    ↪ simple to understand.

⑤ easily overfit if you don't do Laplace smoothing.