← Naive Bayes :→ classification algorithm
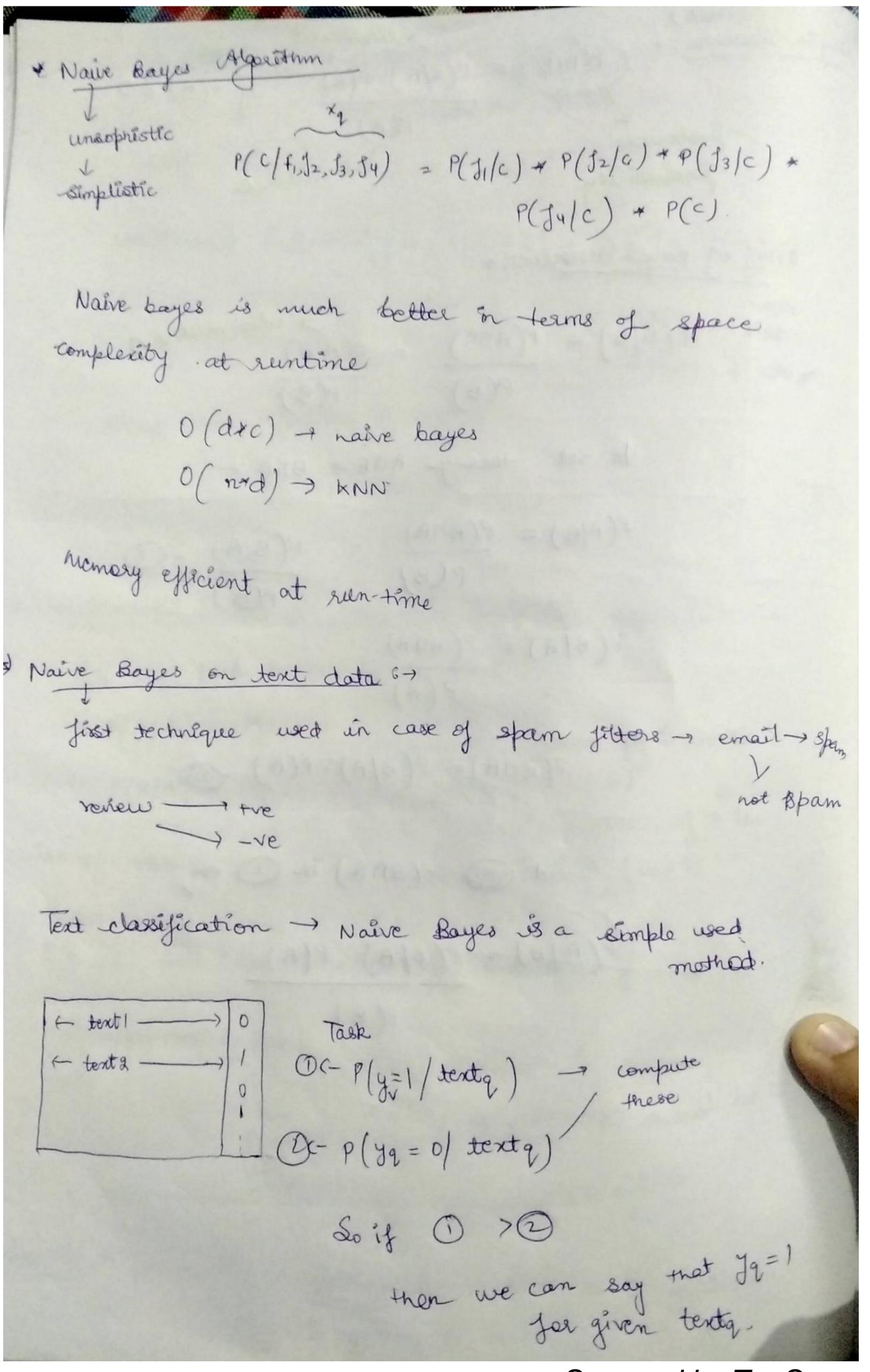
↳ probability - based

KNN→ neighbourhood based classification.

Conditional Probability $= (P(A|B)) = Pr(A = a / B = b)$

↓ 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.

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

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

forex :→

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

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

A, B are said to be mutually exclusive if.

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

$\frac{P(A \cap B)}{P(B)}$     $\frac{P(B \cap A)}{P(A)}$     , so, $P(A \cap B)$ should be 0

as $A \cap B = B \cap B$.

for ex :→ probability of getting 3 in die 1 is 0 if die 2 is getting 6 in it.

**Bayes Theorem :→** (1700's)

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

likelihood → prior → evidence

posterior probability

### Proof of Bayes theorem :→

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

In set theory $A \cap B = B \cap A$ ←

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

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

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

Put ② $P(B \cap A)$ in ① eq

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

Scanned by TapScanner

* Naive Bayes Algorithm

↓
unsophistic
↓
simplistic

$$\underbrace{x_q}$$

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

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

$$O(d*c) \rightarrow \text{naive bayes}$$
$$O(n*d) \rightarrow kNN$$

memory efficient at run-time

② Naive Bayes on text data ↝

first technique used in case of spam filters → email → spam
↓
not spam

review ——→ +ve
——→ -ve

Text classification → Naive Bayes is a simple used method.



Task

① ⟵ $P\left(y=1 / text_q\right)$  → compute these

② ⟵ $P\left(y_q = 0 / text_q\right)$

So if  ① > ②

then we can say that $y_q = 1$ for given $text_q$.

Preprocessing :→      text $\longrightarrow$   stopwords

                                   Stemming

                                   n-grams

$$\hookrightarrow \{ w_1, w_2, w_3, \ldots w_d \}$$

$$\updownarrow$$

                                       Binary BOW

$$\text{text} \xrightarrow{\text{pre-pro}} \underbrace{\{ w_1, w_2, \ldots w_d \}}_{\text{set of words}}$$

$\boxed{P(y_1 = 1 \mid \text{text})} = P(y_i = 1 \mid \underbrace{w_1, w_2, w_3, \ldots, w_d}_{\text{features}})$

↓

*class*
*prior*

                                                     *likelihoods*

$$= P(y=1) * P(w_1 \mid y=1) * \boxed{P(w_2 \mid y=1)}$$

$$\ldots \ldots P(w_d \mid y=1)$$

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

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

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

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

Train data with $y = 1$

Train data with $y = 0$

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

Text - classification problems.

{ spam - detection

Polarity of a review } Naive Bayes is a very good Baseline

↓

benchmark

→ not replacian smoothing

* Laplace - Additive Smoothing :→

After training →

all this data is already computed.

{ $P(y=1)$ ; $P(y=0)$ ← class priors

$P(w_1 | y=1)$    $P(w_1 | y=0)$

$P(w_2 | y=1)$    $P(w_2 | y=0)$

⋮

$P(w_m | y=1)$    $P(w_m | y=0)$ } → likelihd

<u>Test</u> :→ $text_q = (w_1, w_2, w_3, w')$

← *// $w'$ is not present in $\{w_1, w_2, w_3, \ldots \cdots w_m\}$

very often

↗ training data.

$$P(1/text_q) = P(y=1 \mid w_1, w_2, w_3, w')$$

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

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

$P(w'/y=1)$ ⟵

ignoring or dropping it will mean

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

which is not correct.

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

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

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

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

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

Laplace smoothing or additive smoothing :→

$$P(f_1 = a / y = 1) = \frac{0 + \alpha}{n_1 + \alpha k}$$

$\alpha = 1$ → typically (not always)

$k$ = no. of distinct values which $f_1$ can take

$f_1$ ⇒ feature

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

$k = 2$ → because $w'$ is $0$ or $1$.
present or not.

Let $n_1 = 100$

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

*α small, we are getting rid of multiplying all the probabilities with 0.*

$P(w' / y = 1) \neq 0$ ⇓ which implies

$P(y = 1 | textq) \neq 0.$

So, its not 0 anymore

← Case 2 :→ $\alpha = 10000$ → when $\alpha$ is large

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

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

↓

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

## Laplace smoothing :→

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 $n_1$ is small

num or den is small → less confidence in ratio

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

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