

W207.2 Notes 05/22/18

- Review of prob.
- Text-as-data
- Naive Bayes + application

# Probability

A - It will rain today

B - My dog will bark today

→ frequentist

→ Bayesian

A-RAM

$$P(A) = \frac{800}{1000} = \frac{4}{5}$$

$$P(B) = \frac{980}{1000} = \frac{98}{100} \approx 98\%$$

Bayesian - when we calculate probabilities we need to take into account beliefs about probability distribution

$A \perp B$   $A$  - Rain  $B$  - dog bark

$$P(A) = 4/5, \quad P(B) = 9/10$$

$$P(A \text{ and } B) = P(A \cap B) = P(A)P(B) \\ = \left(\frac{4}{5}\right)\left(\frac{9}{10}\right) = \frac{36}{50}$$

$$P(A \text{ or } B) = P(A \cup B) \approx 72\% \\ = P(A) + P(B)$$

$$P(A^c) = 1 - P(A) = 1 - \frac{4}{5} = \frac{1}{5}$$

$$P(B^c) = 1 - P(B) = 1 - \frac{9}{10} = \frac{1}{10}$$

## Bayes Rule

Condition  $\quad$  probability  
likelihood

$$\underbrace{P(A|B)}_{\text{posterior}} = \frac{\underbrace{P(B|A)}_{\text{likelihood}} \underbrace{P(A)}_{\text{prior}}}{\underbrace{P(B)}_{\text{marginal likelihood}}}$$

$$P(B) = P(B|A)P(A) + P(B|A^c)P(A^c)$$

$$P(B) = \underbrace{P(B|A)} \underbrace{P(A)} + \underbrace{P(B|A^c)} \underbrace{P(A^c)}$$

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

$P(A)$   
 $\downarrow$   

$\frac{4}{5}$

$P(\text{Rain} | \text{My dog barked today})$

$$= \frac{(1/2)(4/5)}{(9/10)} = \frac{(4/10)}{9/10} = \frac{4}{9} \approx 44\%$$

Assn. For Bayes Thm  $A \perp B$

$$P(A) = 9/10$$

$$P(A|B) = \frac{(1/2)(9/10)}{(9/10)} = 50\%$$



# Text-as-DatA

$N = 100$  email (documents)

(Corpus)

Spam  
Target/  $\begin{cases} 1 & \text{email is spam} \\ 0 & \text{it is not spam} \end{cases}$

(class label)

(Features)

Text - words/terms

= "you have won one million dollars give bank info"



Spam

1

0

0

1

.

1

.

1

" Text  
you have win . - "

$$\boxed{\text{Spam}} = f(\text{Text})$$
$$\hat{\text{Spam}} = \hat{f}(\text{Text})$$

# Text to Data

## ① Preprocessing

→ tokenization

→ cleaning

- o stemming

- o punctuation removal

- o stop word removal

NLP

## ② Building a DTM (Document term matrix)

# Pre-processing

## ① Tokenization

"you have win one million..."

a) partition text into terms

n-grams    unigrams (words)

$n \geq 2$  - phrase    2-grams (phrase)

---

$$W_1 = \{ \text{"you", "have", "win", "one", "million", "you have" ...} \}$$

$w_i \in \mathbb{R}^g$  (b) Clean texts

— Stemming - groups  $\rightarrow$  group

groupings  
grouped  $\rightarrow$  group

won  $\rightarrow$  win

wins  $\rightarrow$  win

-ing  
-ed  
-s

$\rightarrow$  remove punctuation, remove it's  
lowercase

$\rightarrow$  Stop word removal - removal of common and frequent terms

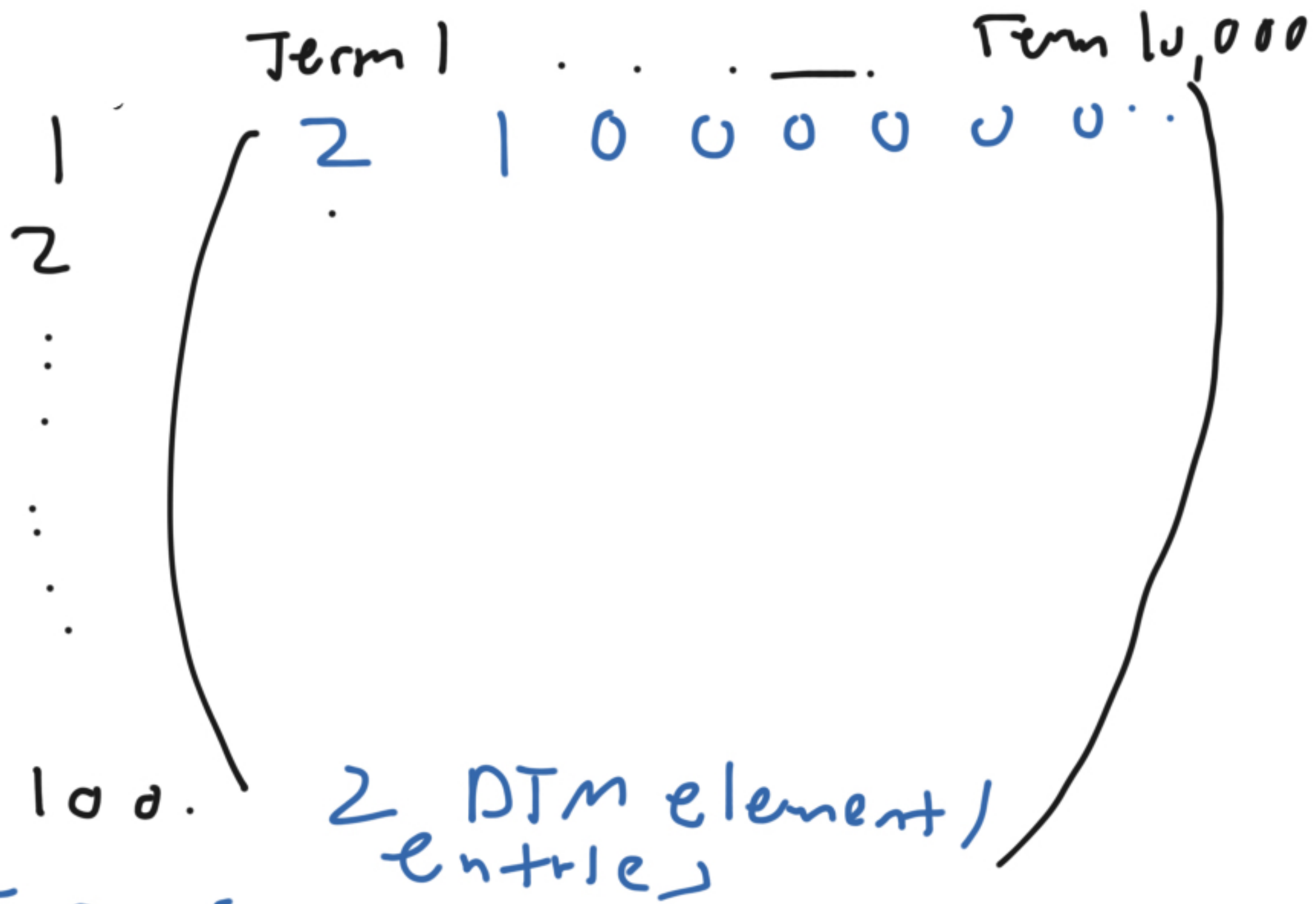
$w_1 = \{ \text{"you", "have", "when", ...} \}$

② DTM

$X \in \mathbb{R}^{N \times W}$   $N \equiv \text{obs / documents}$   
 $W \equiv \text{terms in all}$   
 $N \text{ documents}$

$N = 100 \text{ docs / email}$

$W = 10,000 \text{ terms}$   $X \in \mathbb{R}^{100 \times 10000}$



① Term frequency



$w_1 = \{ \text{'you', 'have', 'win' ...} \}$

<u>Doc</u>	<u>Span</u>	<u>you</u>	<u>have</u>	<u>win</u>	<u>.</u>	<u>-</u>
1	1	1	1	1	0	0
2	0	0	0	0		

② TF-IDF - Term-frequency  
 Inverse Document frequency

you have wm

$1(\frac{1}{80})$   $1(\frac{1}{50})$   $1(\frac{1}{5})$

$IDF = \frac{1}{\text{\# of documents word/term } i \text{ appears in}}$

$$IDF(you) = \frac{1}{80} \quad IDF(have) = \frac{1}{50}$$

$$IDF(wm) = \frac{1}{5}$$

$X \equiv DTM$  or Feature matrix

$N \times W$

You have nm...

$N=2100$

$\begin{pmatrix} 1 \\ \cdot \\ \cdot \\ : \\ 100 \end{pmatrix}$

$$\text{Target} = \text{spam} \in \mathbb{R}^{100}$$
$$X \in \mathbb{R}^{100 \times 10000}$$

Naive Bayes

$$N=100$$

- ① Pre-processing / DTM
- ② Divide data into training and test 80% / 10%
- ③ Train algorithm

④ Monitor performance  
80 training emails

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

$$P(S=1|w_i) = \frac{P(w_i|S=1) P(S=1)}{P(w_i)}$$

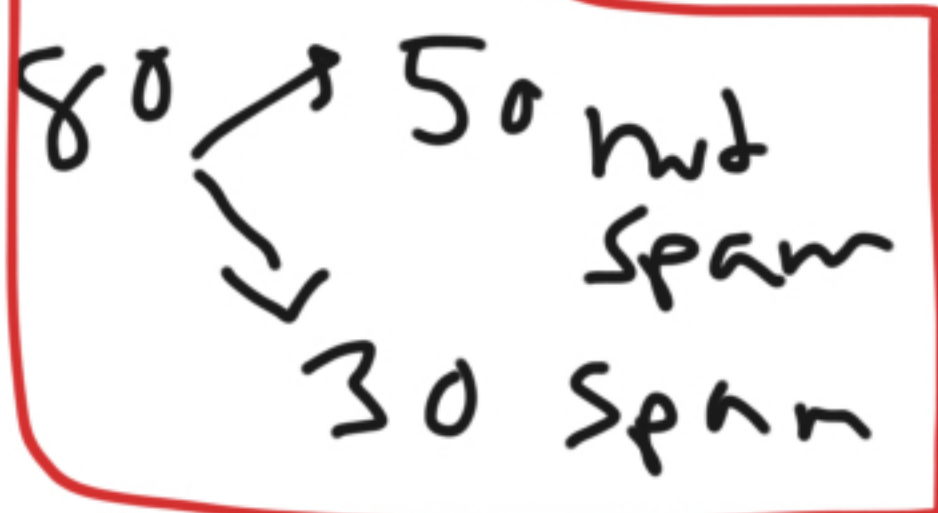
Likelihood

$w_1 = \{ \text{"you"}, \text{"have"}, \text{"win"} \}$

$$P(w_1 | S=1) = \frac{P(\text{"you"}, \text{"have"}, \text{"win"} | S=1)}{1/30}$$

$$\Rightarrow P(\text{"you"} | S=1) P(\text{"have"} | S=1)$$

$$P(\text{"win"} | S=1) = 15/30$$



$$P(\text{"you"} | S=1) = \frac{5}{30}$$



$$P(w_1 | S=1) = \left(\frac{5}{30}\right) \left(\frac{1}{30}\right) \left(\frac{15}{30}\right)$$

$$P(S=1) = \frac{30}{50} = \frac{3}{5} \quad P(S=0) = \frac{20}{50}$$

$$P(w_1) = P(w_1 | S=1)P(S=1) + P(w_1 | S=0)P(S=0)$$

$$P(S=1 | w_1) = \frac{\binom{75}{27000} \left(\frac{3}{5}\right)}{\binom{75}{27000} \left(\frac{3}{5}\right) + \binom{20}{27000} \left(\frac{2}{5}\right)}$$

||