

AI VIETNAM
All-in-One Course
(TA Session)

Text Classification

Project



AI VIET NAM
[@aivietnam.edu.vn](http://aivietnam.edu.vn)

Dinh-Thang Duong – TA
Anh-Khoi Nguyen – STA

Getting Started

◆ Objectives

$$p(Y|X) = \frac{p(X|Y) \times p(Y)}{p(X)}$$

Spam	Not spam
 <p>"SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply Reply HL 4 info"</p>	 <p>"Nah I don't think he goes to usf, he lives around here though"</p>

Our objectives:

- Discuss about Naïve Bayes algorithm.
- Delve into one of popular NLP tasks: Text Classification.
- Apply Naïve Bayes and its variants to solve a text classification task.
- Investigate an improved baseline over Naïve Bayes and conduct experiments.

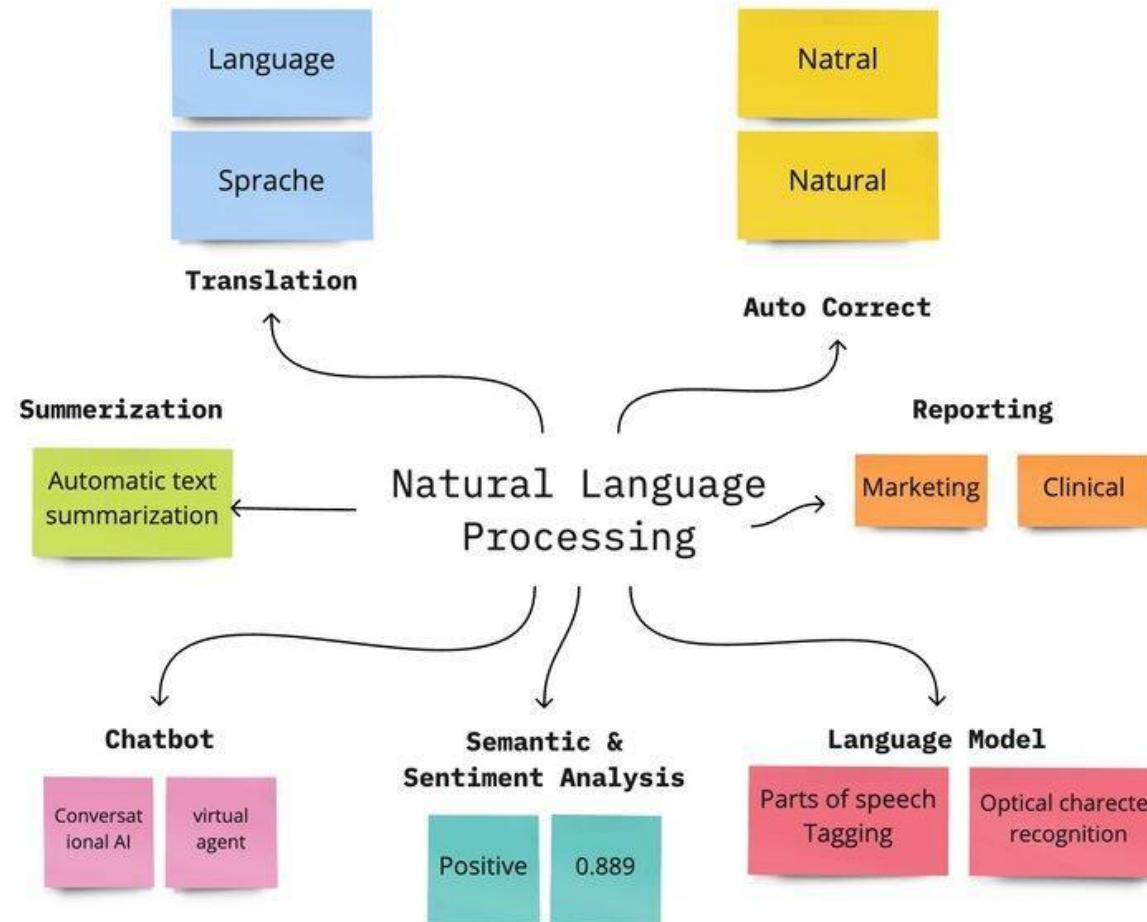
Outline

- Introduction
- Message Classification
- Code Implementation
- Improved Baseline
- Question

Introduction

Introduction

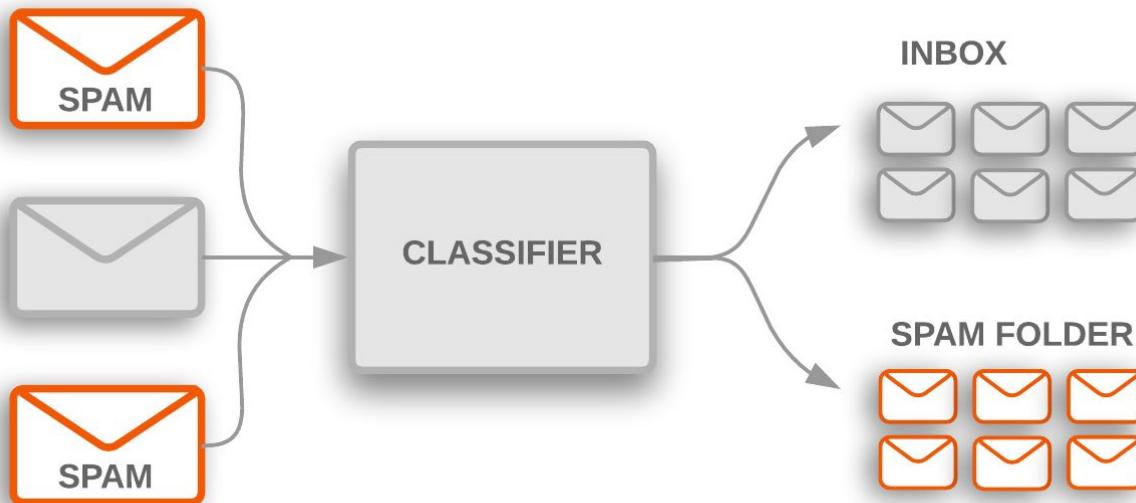
◆ Getting Started



miro

Introduction

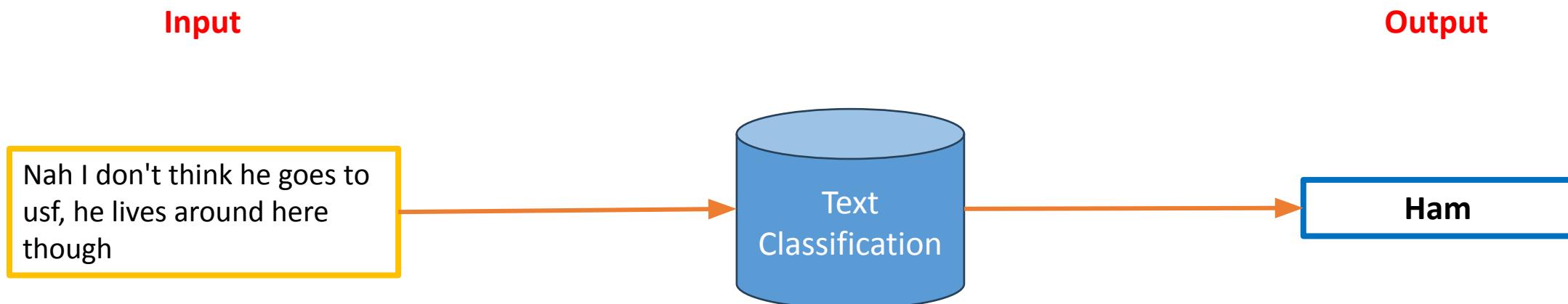
❖ Text Classification



Text classification: A Natural Language Processing (NLP) task that involves categorizing text into predefined labels or classes. It is used to automatically assign a category to a text document, such as spam detection in emails, sentiment analysis of reviews, or topic classification of articles.

Introduction

❖ Text Classification I/O



Text Classification: An NLP task that aims to classify a given text into pre-defined classes.

Introduction

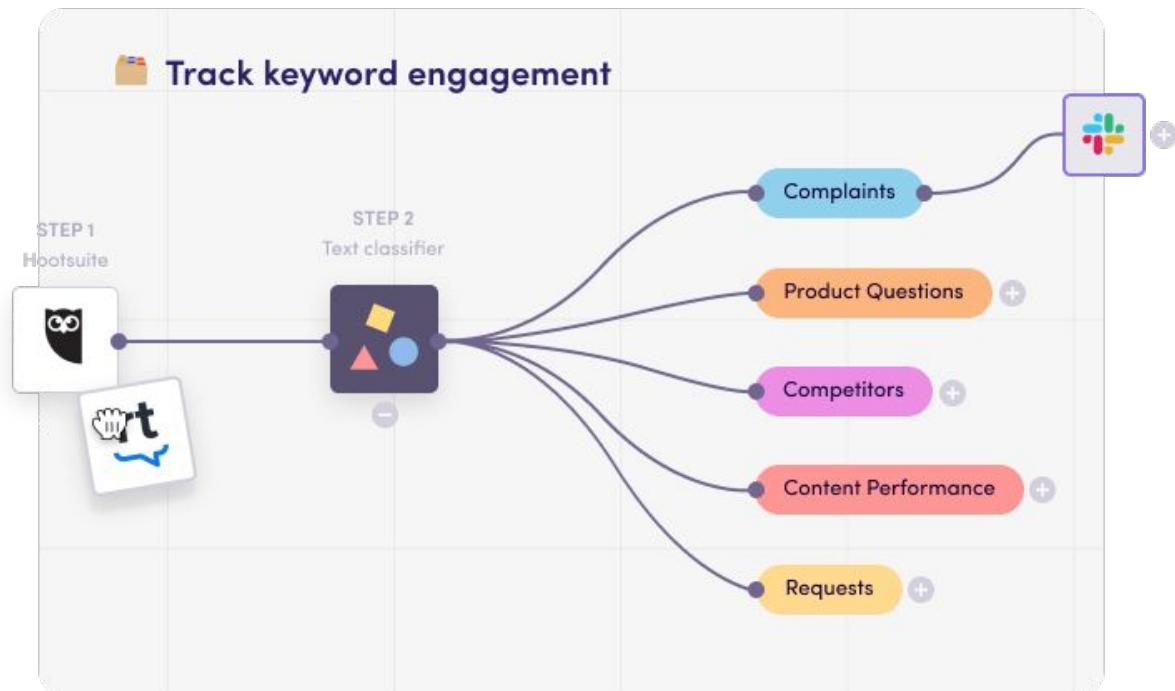
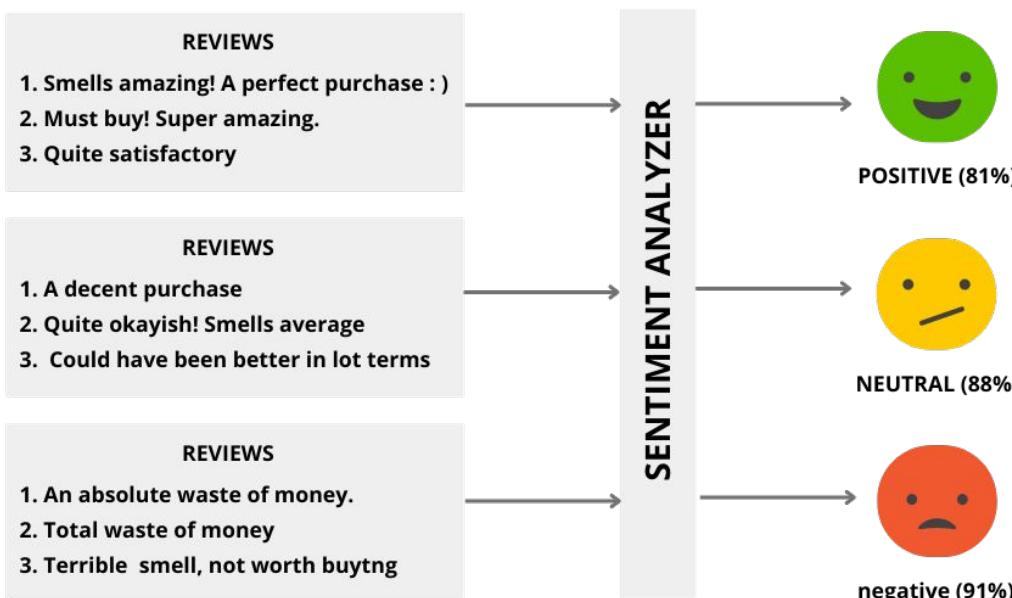
❖ Applications



Fragrance-1
(Lavender)

Fragrance-1
(Rose)

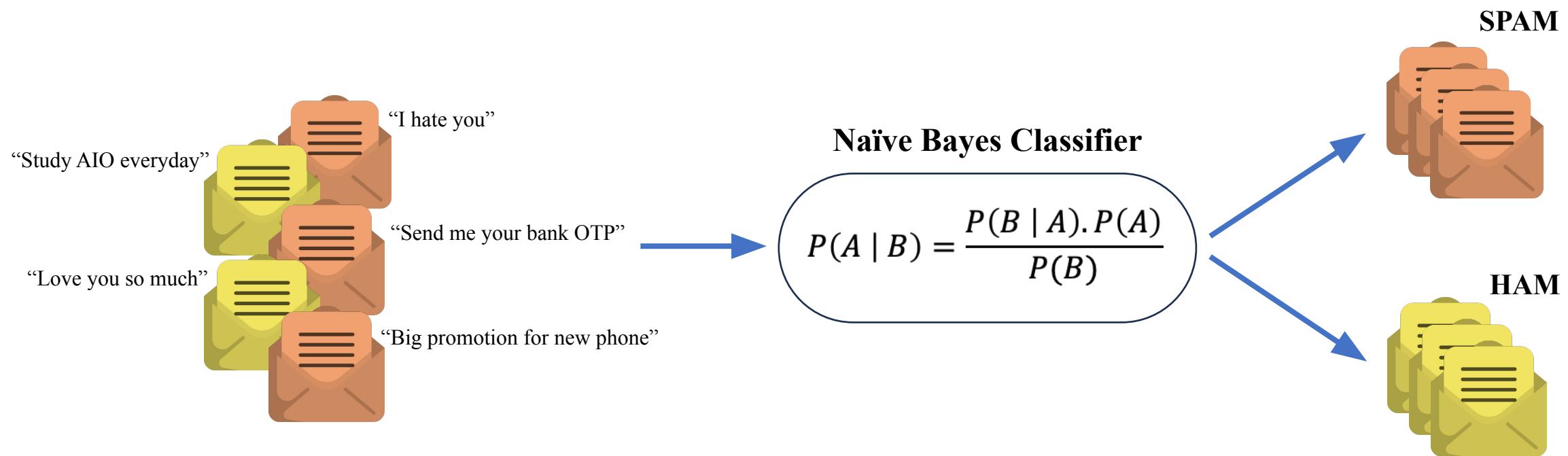
Fragrance-1
(Lemon)



Introduction

◆ Project Statement: Text Classification Naïve Bayes

Description: Given Message Classification Dataset, build a Naives Bayes model to determine whether a text message is spam message or not (ham).



Message Classification

Message Classification

◆ Probability

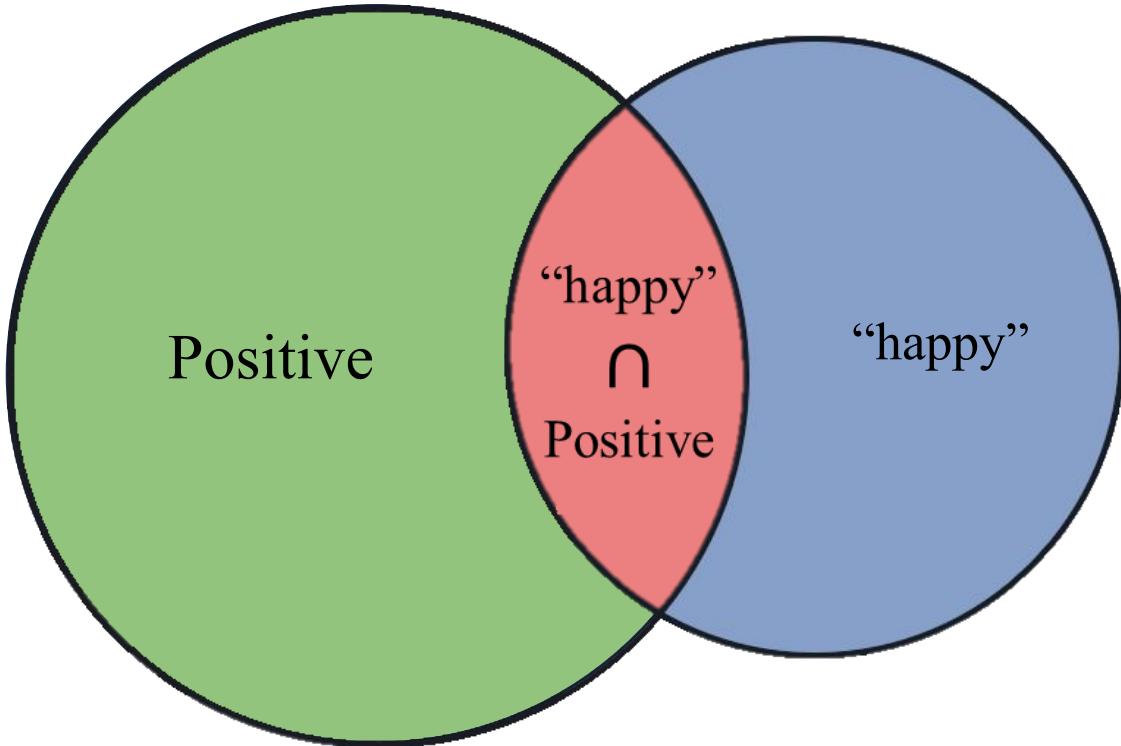
Positive	Positive	Positive	Positive	Positive
Positive	Positive	Positive	<i>Negative</i>	<i>Negative</i>
<i>Negative</i>	<i>Negative</i>	<i>Negative</i>	<i>Negative</i>	<i>Negative</i>
<i>Negative</i>	<i>Negative</i>	<i>Negative</i>	<i>Negative</i>	<i>Negative</i>

$$P(Pos) = \frac{N_{Pos}}{N} = \frac{8}{20} = 0.4$$

$$P(Neg) = 1 - P(Pos) = 1 - 0.4 = 0.6$$

Message Classification

◆ Conditional Probability



$$P("happy" | Pos) = \frac{P(Pos \cap "happy")}{P(Pos)}$$

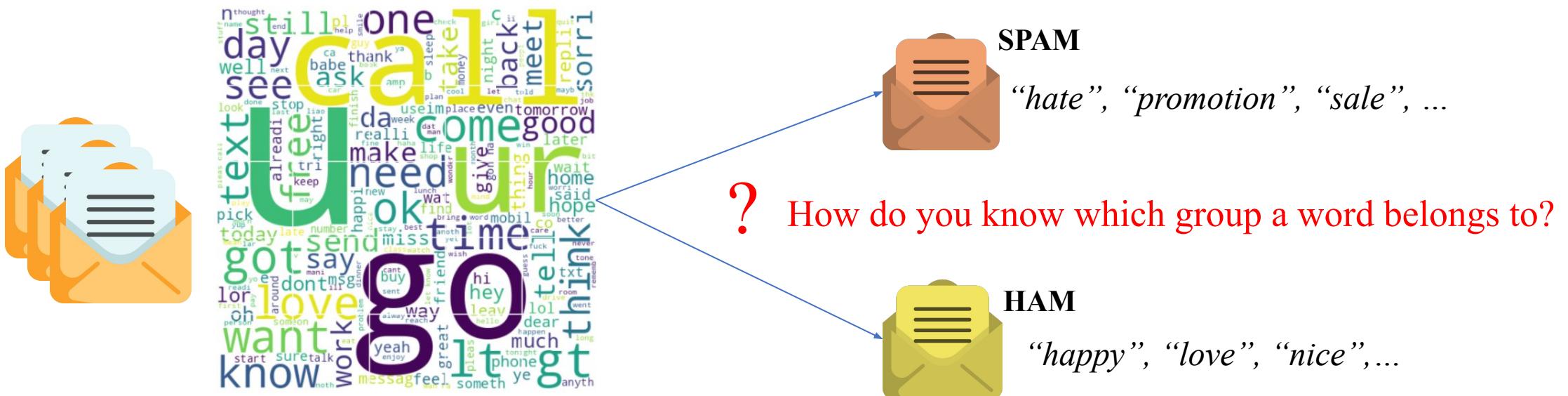
Message Classification

◆ Bayes' Rule

The content of a letter consists of many words combined together

How do we know which letter is SPAM or HAM?

One of the ways is based on individual words



Message Classification

◆ Bayes' Rule

How do you know which group a word belongs to?



Bayes' Rule: $P(A | B) = \frac{P(B | A).P(A)}{P(B)}$

$$P(\text{Class} | \text{Words}) = \frac{P(\text{Words} | \text{Class}) . P(\text{Class})}{P(\text{Words})}$$

$$P(\text{Ham} | w_1, \dots, w_n) = \frac{P(w_1, \dots, w_n | \text{Ham}).P(\text{Ham})}{P(w_1, \dots, w_n)}$$

!!!

$$P(\text{Ham} | "happy", "hate", "nice", \dots) = \frac{P("happy", "hate", "nice", \dots | \text{Ham}).P(\text{Ham})}{P("happy", "hate", "nice", \dots)}$$

Message Classification

◆ Word Dependence

$$P("brown", "sugar", "pearl", "milk", "tea" | \text{Ham})$$

Human:



“brown” sugar

pearl milk tea”

Naïve:



$$P("brown" | \text{Ham}).P("sugar" | \text{Ham}).P("pearl" | \text{Ham}).P("milk" | \text{Ham}).P("tea" | \text{Ham})$$

Message Classification

❖ Naïve Bayes

Bayes' rule with the “Naïve” is the presence of one feature (word) does not affect the presence of another (word).

$$P("happy", "hate", "nice", \dots | Ham) = \prod_{i=1}^n P(word_i | y)$$

$= D("happy" | Ham) \cdot D("hate" | Ham) \cdot D("nice" | Ham) \cdot D(\dots | Ham)$

$$P(Ham | "happy", "hate", "nice", \dots) = \frac{\prod_{i=1}^n P(word_i | Ham) \cdot P(Ham)}{P("happy", "hate", "nice", \dots)}$$

Message Classification

❖ Remove constant

This value does not change when we consider different y classes.

$$P(\text{Ham} | "happy", "hate", "nice", \dots) = \frac{\prod_{i=1}^n P(\text{word}_i | \text{Ham}) \cdot P(\text{Ham})}{P("happy", "hate", "nice", \dots)}$$

$$P(\text{Spam} | "happy", "hate", "nice", \dots) = \frac{\prod_{i=1}^n P(\text{word}_i | \text{Spam}) \cdot P(\text{Spam})}{P("happy", "hate", "nice", \dots)}$$

Message Classification

❖ Remove constant

This value does not change when we consider different y classes.

$$P(\text{Ham} | "happy", "hate", "nice", \dots) = \frac{\prod_{i=1}^n P(\text{word}_i | \text{Ham}) \cdot P(\text{Ham})}{P("happy", "hate", "nice", \dots)}$$

$$P(\text{Spam} | "happy", "hate", "nice", \dots) = \frac{\prod_{i=1}^n P(\text{word}_i | \text{Spam}) \cdot P(\text{Spam})}{P("happy", "hate", "nice", \dots)}$$

Naïve Bayes Mail Classification :

Message Classification

❖ Project Statement

Description: Given Message Classification Dataset, build a Naives Bayes model to determine whether a text message is spam message or not (ham).

Spam	Not spam
 <p>“SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply Reply HL 4 info”</p>	 <p>“Nah I don't think he goes to usf, he lives around here though”</p>

Message Classification

❖ Introduction

Email classification is a task where we determine whether a given Mail is

- Spam (unsolicited or unwanted) Spam
- Ham (non-spam) Ham

Suppose you are “all in” AIO and you receive the following emails:

 Chính	 Qu... 49 cuộc trò chuyện mới DeepLearning.AI, Ivan at Notion...	 Mạn... 11 cuộc trò chuyện mới LinkedIn, YouTube	 Nội... 50 cuộc trò chuyện mới LinkedIn Job Alerts, ngrok, Real ...
<input type="checkbox"/> ★ AI VIET NAM	Ham	M02EC05 - SEMINAR Confirmation - Hello, Thank you for registering for M02EC05 - SEMINAR. You can find information about...	6 thg 8
<input type="checkbox"/> ★ Scammer	Spam	✨ Send the bank OTP - To receive the special gift, you need to send me the OTP code	6 thg 8
<input type="checkbox"/> ★ Friend	Spam	🍴 Invite to eat hot pot - Take a break from AIO class today to go eat, today Haidilao has a promotion	6 thg 8
<input type="checkbox"/> ★ AI VIET NAM	Ham	M02EC03 - Basic Statistics Confirmation - Hello, Thank you for registering for M02EC03 - Basic Statistics. You can find infor...	6 thg 8
<input type="checkbox"/> ★ Ads	Spam	💻 Laptop promotion - Super promotions on laptops and phones, buy now	6 thg 8

Message Classification

◆ Message Classification I/O

Spam	Not spam
 <p>"SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply Reply HL 4 info"</p>	 <p>"Nah I don't think he goes to usf, he lives around here though"</p>

Problem Statement: Given a message, classify it into one of the two classes: Spam or Ham

Input: "Nah I don't think he goes to usf, he lives around here though"

Output: "Ham"

Message Classification Problem

Message Classification

❖ Dataset

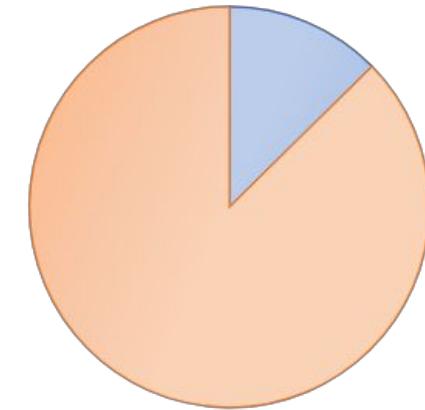
The data set includes 2 columns:

- Category (mail type)
- Message (mail content)

▲ Category	▲ Message
ham	87%
spam	13%
	5157 unique values
ham	Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got a...
ham	Ok lar... Joking wif u oni...
spam	Free entry in 2 a wkly comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entr...
ham	U dun say so early hor... U c already then say...
ham	Nah I don't think he goes to usf, he lives around here though
spam	FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like some fun you up for

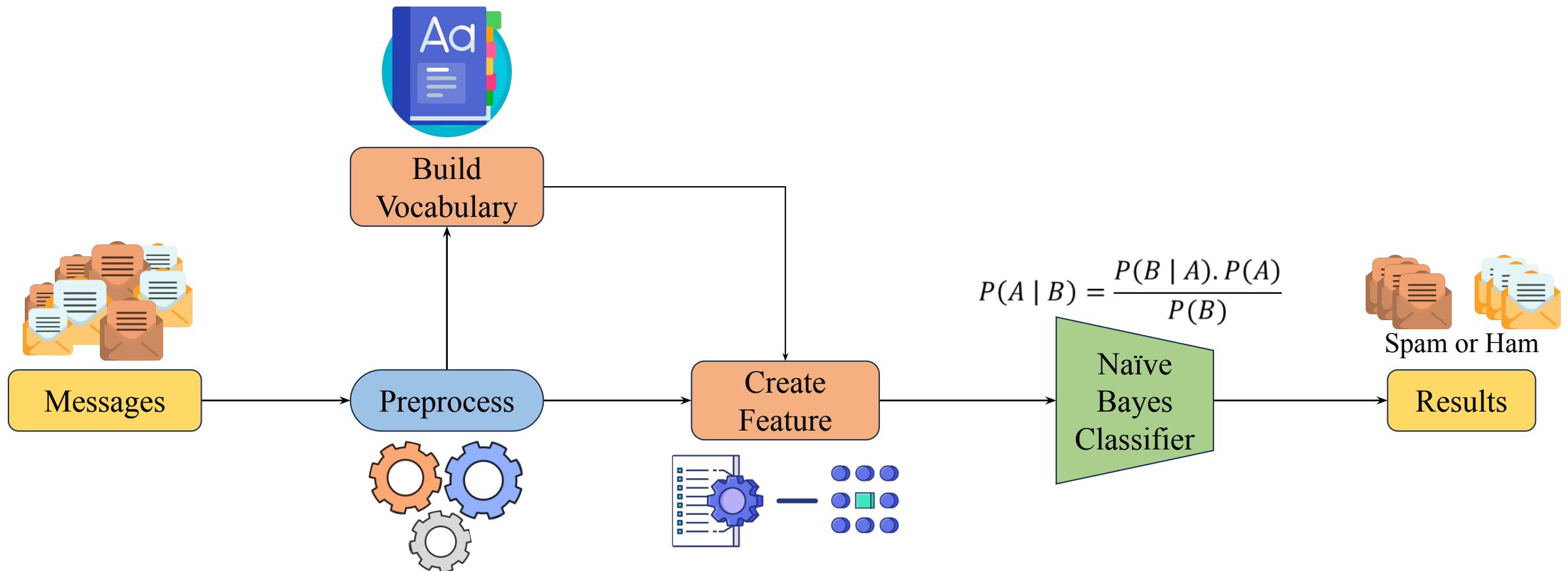
Email Classification

■ Spam ■ Ham



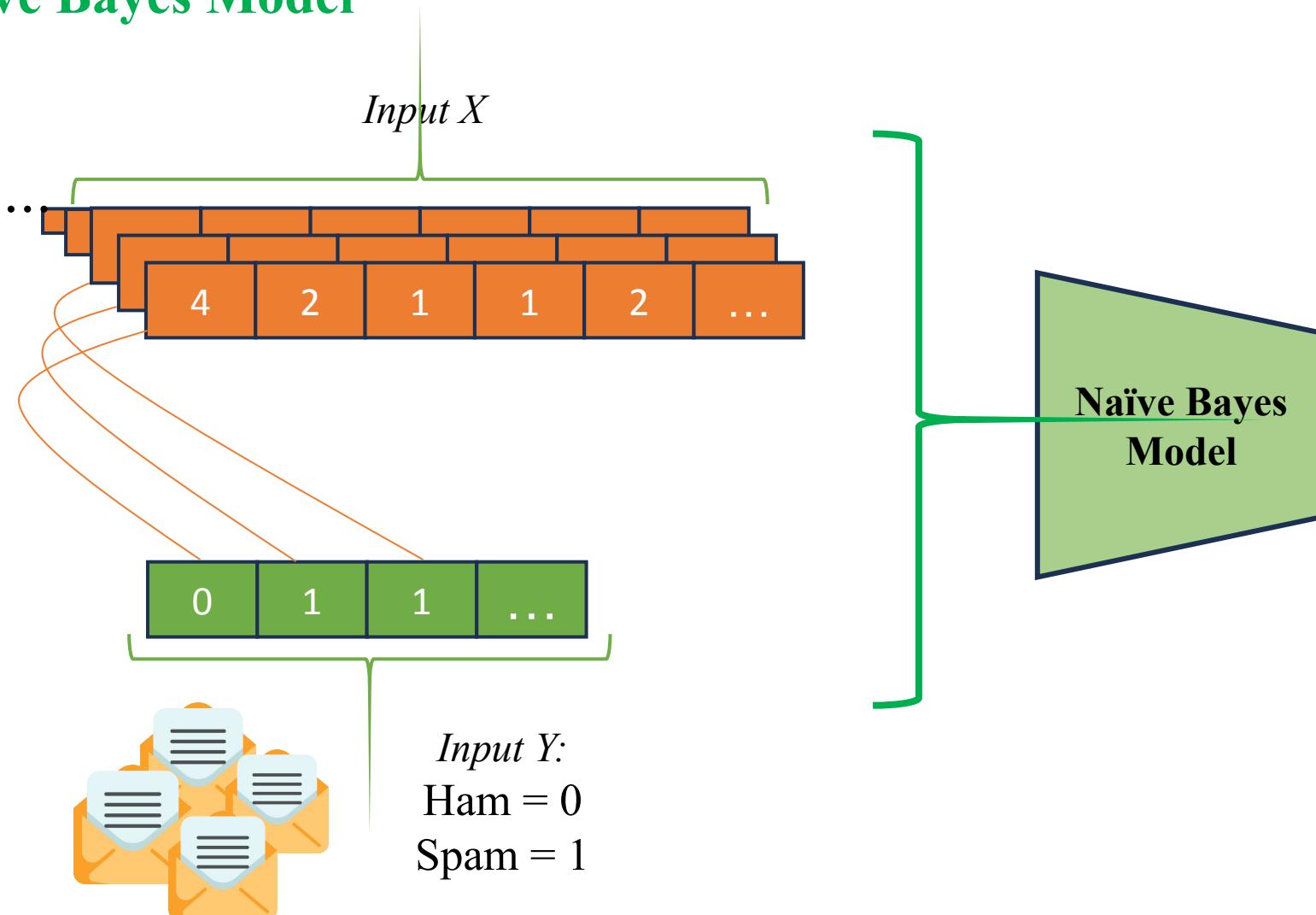
Message Classification

◆ Project Pipeline



Message Classification

◆ Input Naïve Bayes Model

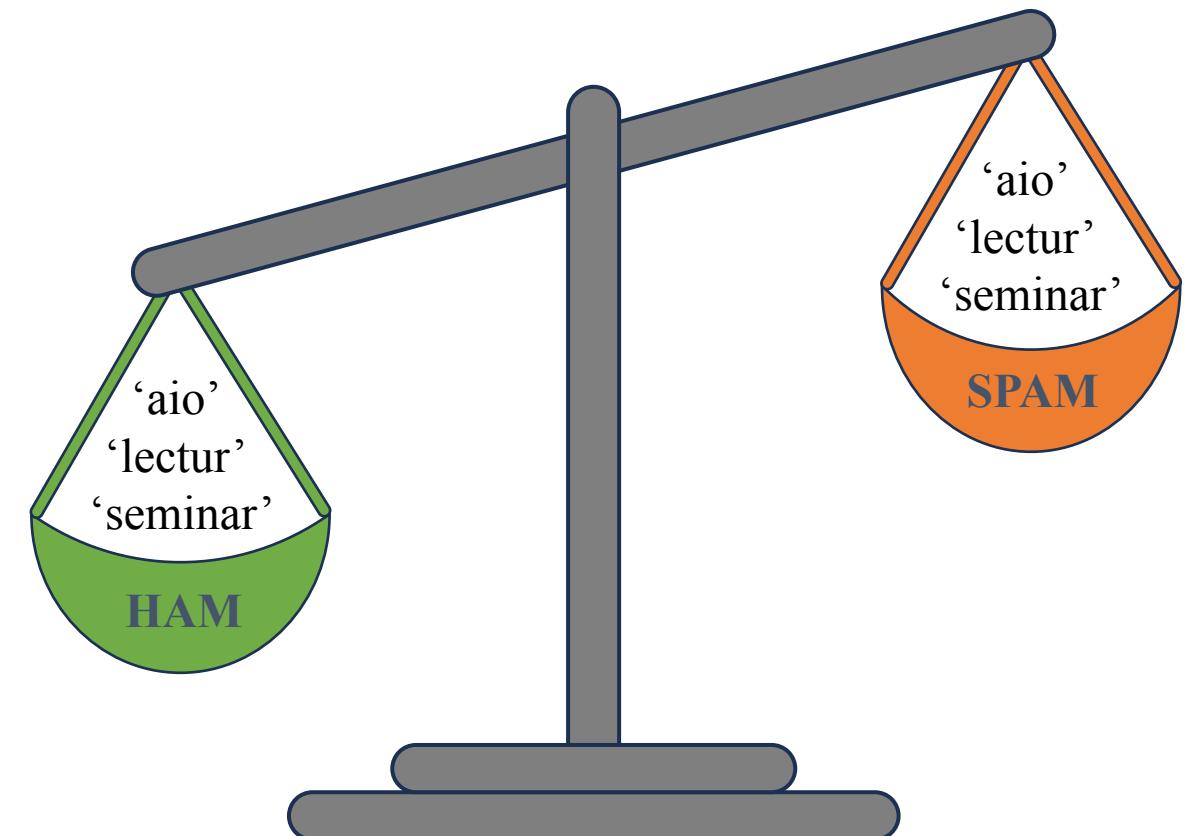


Message Classification

◆ Applying Naïve Bayes with Maximum A Posteriori Estimation

To classification, using MAP to estimation $P(y)$ and $P(word_i | y)$, where y is a label in training set.
Which y has the highest probability is the model result

$$\hat{y} = \arg \max_y P(y) \cdot \prod_{i=1}^n P(word_i | y)$$



Message Classification

❖ $P(x_i|y)$



$$\hat{y} = \arg \max_y P(y) \cdot \prod_{i=1}^n P(x_i | y)$$

AIO	I	study	we	love
-----	---	-------	----	------

4	2	1	1	2
---	---	---	---	---

Suppose Input X:

0	1	1	0	1
---	---	---	---	---

0	3	0	2	0
---	---	---	---	---

Suppose Input Y:

0	0	1
---	---	---

Word	Count	
Word	Ham	Spam
aio	$4 + 0 = 4$	0
i	$2 + 1 = 3$	3
studi	$1 + 1 = 2$	0
we	$1 + 0 = 1$	2
love	$2 + 1 = 3$	0
V	$4+3+2+1+3=13$	$3+2=5$

Message Classification

❖ $P(x_i|y)$



$$\hat{y} = \arg \max_y P(y) \cdot \prod_{i=1}^n P(x_i | y)$$

Word	Ham	Spam
aio	4 / 13	0 / 5
i	3 / 13	3 / 5
studi	2 / 13	0 / 5
we	1 / 13	2 / 5
love	3 / 13	0 / 5



Word	Ham	Spam
aio	0.308	0
i	0.231	0.6
studi	0.154	0
we	0.077	0.4
love	0.231	0

V 13 5

V 13 5

Message Classification

❖ $P(y)$

*Naïve Bayes
Model*

$$\hat{y} = \arg \max_y \mathbf{P}(y) \cdot \prod_{i=1}^n P(x_i | y)$$

Suppose Input Y:

0	0	1
---	---	---

$$P(Ham) = \frac{2}{N} = \frac{2}{3} = 0.667$$

$$P(Spam) = \frac{1}{N} = \frac{1}{3} = 0.334$$

Message Classification

◆ Predict

Naïve Bayes Model

Word	Ham	Spam	$P(Ham) = 0.667$
aio	0.308	0	$P(Spam) = 0.334$
i	0.231	0.6	
studi	0.154	0	
we	0.077	0.4	
love	0.231	0	

V 13 5

Predict

New mail: “I study AIO, I love it”

$$\hat{y} = \arg \max_y P(y) \cdot \prod_{i=1}^n P(word_i | y)$$

$$\begin{aligned}\hat{y}_{Ham} &= (0.667) \cdot (0.231 * 0.154 * 0.308 * 0.231 * 0.231) \\ &\approx 0.000389971807067304\end{aligned}$$

$$\begin{aligned}\hat{y}_{Spam} &= (0.334) \cdot (0.6 * 0.01 * 0.01 * 0.6 * 0.01) \\ &\approx 0.00000012024\end{aligned}$$

$$\hat{y}_{Ham} > \hat{y}_{Spam} \rightarrow Predict = Ham$$

Message Classification

◆ Problem 1 – Zero

Naïve Bayes Model

$$\hat{y} = \arg \max_y P(y) \cdot \prod_{i=1}^n P(\text{word}_i | y)$$

Word	Ham	Spam
aio	4 / 13	0 / 5
i	3 / 13	3 / 5
studi	2 / 13	0 / 5
we	1 / 13	2 / 5
love	3 / 13	0 / 5

Word	Ham	Spam
aio	0.308	0 !
i	0.231	0.6
studi	0.154	0 !
we	0.077	0.4
love	0.231	0 !

V

13

5

V

13

5

Message Classification

◆ Problem 1 - Laplacian Smoothing

Word	Ham	Spam		Word	Ham	Spam
aio	4	0		aio	$(4+1)/(13+5)$	$(0+1)/(5+5)$
i	3	3	$P(x_i y) = \frac{P(x_i) + 1}{V + N_{dictionary}}$	i	$(3+1)/(13+5)$	$(3+1)/(5+5)$
studi	2	0		studi	$(2+1)/(13+5)$	$(0+1)/(5+5)$
we	1	2		we	$(1+1)/(13+5)$	$(2+1)/(5+5)$
love	3	0		love	$(3+1)/(13+5)$	$(0+1)/(5+5)$
V	13	5		V	13	5

Message Classification

◆ Problem 1 - Laplacian Smoothing

$$P(x_i | y) = \frac{P(x_i) + 1}{V + N_{dictionary}}$$

Word	Ham	Spam	Word	Ham	Spam
aio	$(4+1)/(13+5)$	$(0+1)/(5+5)$	aio	0.278	0.1
i	$(3+1)/(13+5)$	$(3+1)/(5+5)$	i	0.222	0.4
studi	$(2+1)/(13+5)$	$(0+1)/(5+5)$	studi	0.167	0.1
we	$(1+1)/(13+5)$	$(2+1)/(5+5)$	we	0.111	0.3
love	$(3+1)/(13+5)$	$(0+1)/(5+5)$	love	0.222	0.1
			Sum	1	1

Message Classification

◆ Problem 1 - Laplacian Smoothing

Predict

Word	Ham	Spam	New mail: "I study AIO, I love it."
aio	0.278	0.1	$\hat{y}_{Ham} = (0.667) \cdot (0.222 * 0.167 * 0.278 * 0.222 * 0.222)$ ≈ 0.00050794
i	0.222	0.4	
studi	0.167	0.1	
we	0.111	0.3	$\hat{y}_{Spam} = (0.334) \cdot (0.1 * 0.4 * 0.1 * 0.3 * 0.1)$ ≈ 0.00012
love	0.222	0.1	
Sum	1	1	$\hat{y}_{Ham} > \hat{y}_{Spam} \rightarrow Predict = Ham$

Message Classification

◆ Problem 2 – Log Likelihood

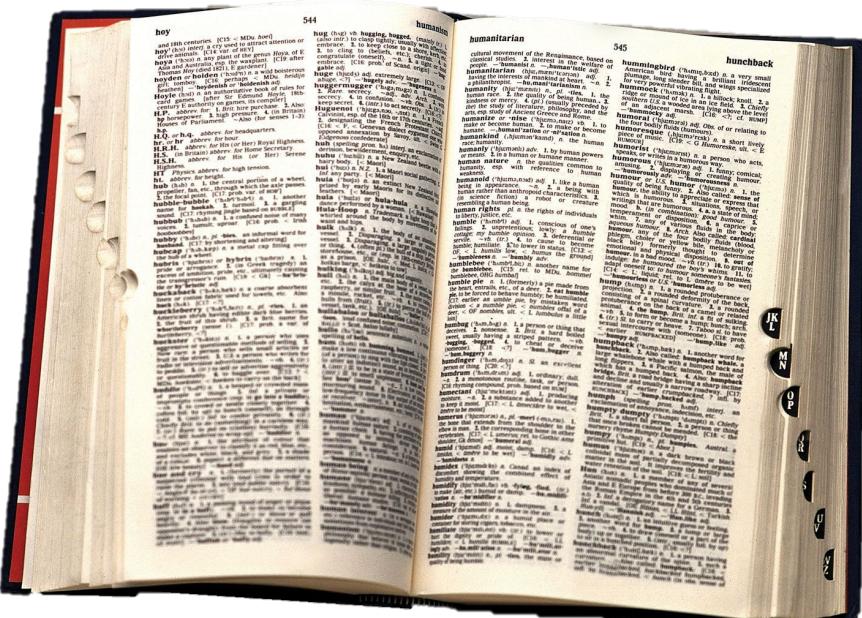
$$\hat{y}_{Ham} \approx 0.00050794$$

!

$$\hat{y}_{Spam} \approx 0.00012$$

!

In reality, the dictionary is huge



The results will become very very small, gradually approaching 0

and the message content much longer

Black Friday in July is here!  Thú rác x



IIN



We couldn't agree more. To kick off our Black Friday in July sale, this is why we think summer is the best time to enroll in one of our courses:

→ Take the NBHWC board-certification exam by 2025. You can save 30% on tuition for the Health Coach Board Certification Training, which will prepare you to take the NBHWC (National Board for Health & Wellness Coaching) board-certification exam by next year! [Click here and use code BFINJULY to get started.](#)

→ The Health Coach Training Program is about to start. Transform your life, launch your coaching career, and take 30% off! Class starts July 22nd, so don't wait. [Click here and use code BFINJULY to get started.](#)

→ Black Friday deals ... in July! Summer can only be made better with an early Black Friday sale! Take 30% off not just the courses above, but ALL COURSES for a limited time only.

Message Classification

◆ Problem 2 – Log Likelihood

Take advantage of the **Scaling** property of *Logarithms*.

0.00001	-5
0.001	-3
0	Error
100	3
10000	5

$$\prod_{i=1}^n P(\text{word}_i \mid y) = \sum_{i=1}^n \log(P(\text{word}_i \mid y))$$

Message Classification

◆ Problem 2 – Log Likelihood

Predict

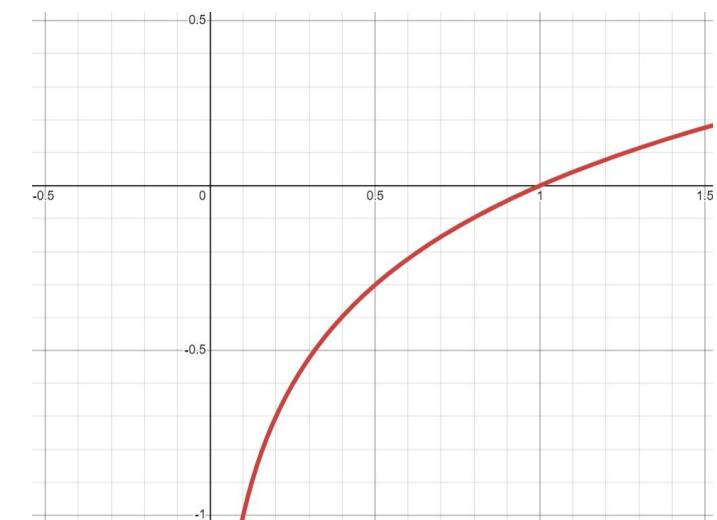
New mail: “**I study AIO, I love it.**”

$$\begin{aligned}\hat{y}_{Ham} &= (0.667) \cdot (\log(0.222) + \log(0.167) + \log(0.278) + \log(0.222) + \log(0.222)) \\ &\approx (0.667) \cdot (-3.294) \approx -2.197098\end{aligned}$$

$$\begin{aligned}\hat{y}_{Spam} &= (0.334) \cdot (\log(0.1) + \log(0.4) + \log(0.1) + \log(0.3) + \log(0.1)) \\ &\approx (0.334) \cdot (-3.921) \approx -1.30955346\end{aligned}$$

Because taking the Logarithms gives a value less than 1, the label

$$\hat{y}_{Ham} < \hat{y}_{Spam} \rightarrow Predict = Ham$$

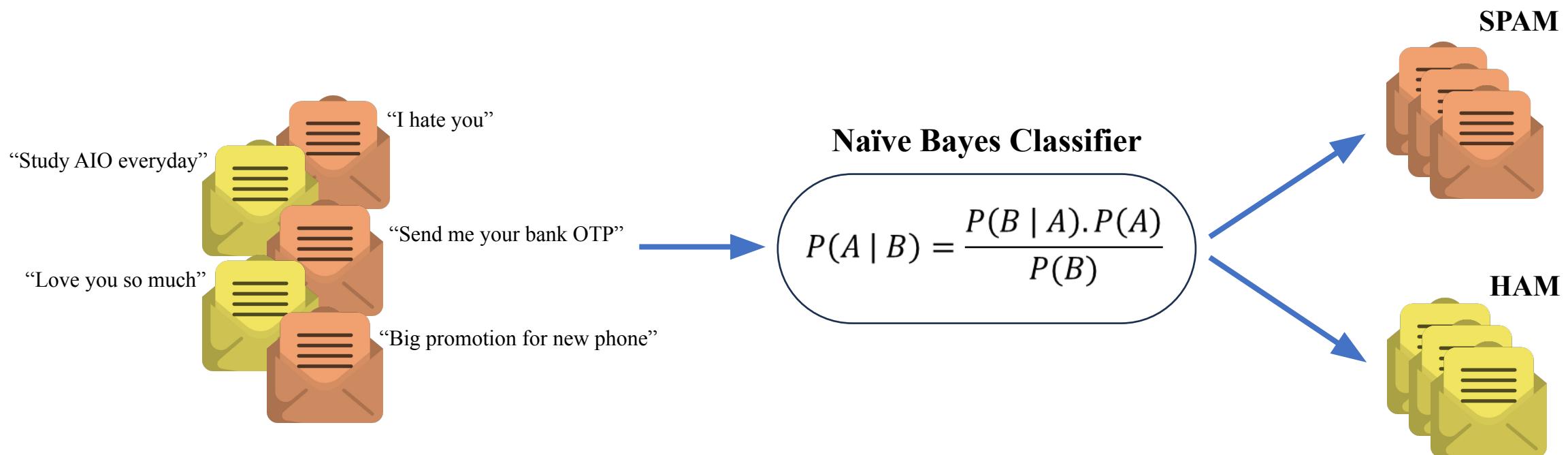


Code Implementation

Code Implementation

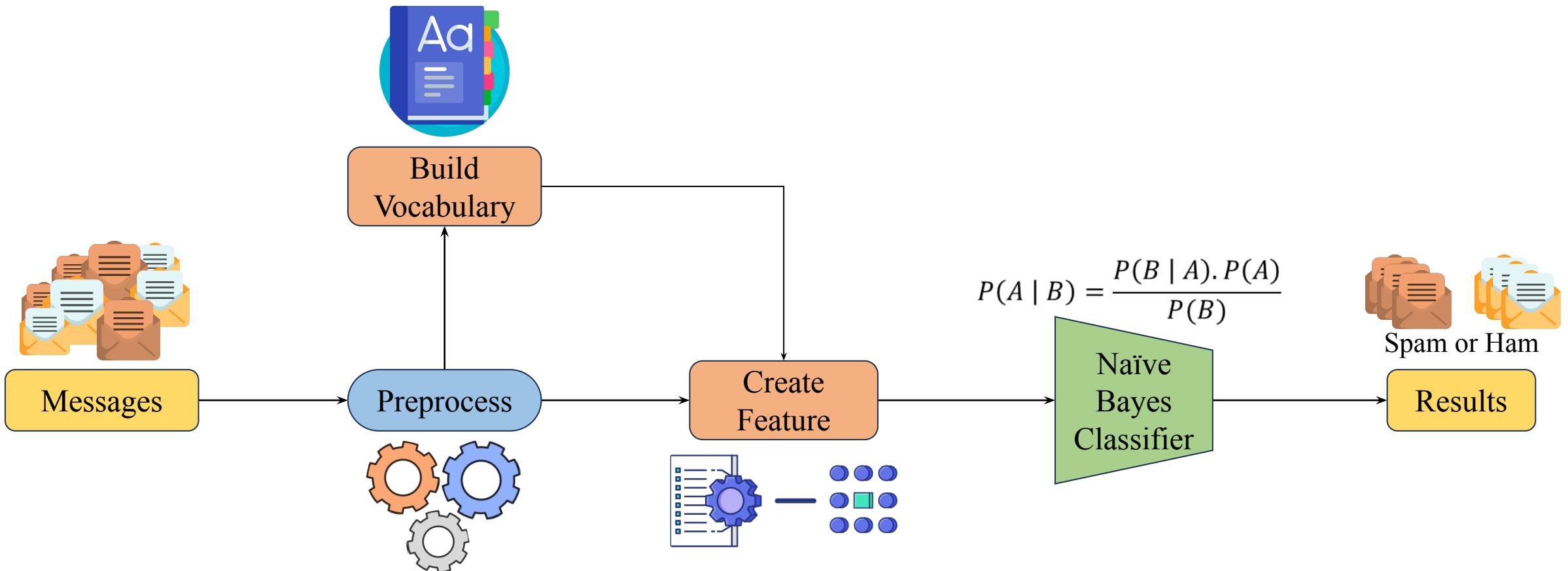
❖ Introduction

Description: Given Message Classification Dataset, build a Naives Bayes model to determine whether a text message is spam message or not (ham).



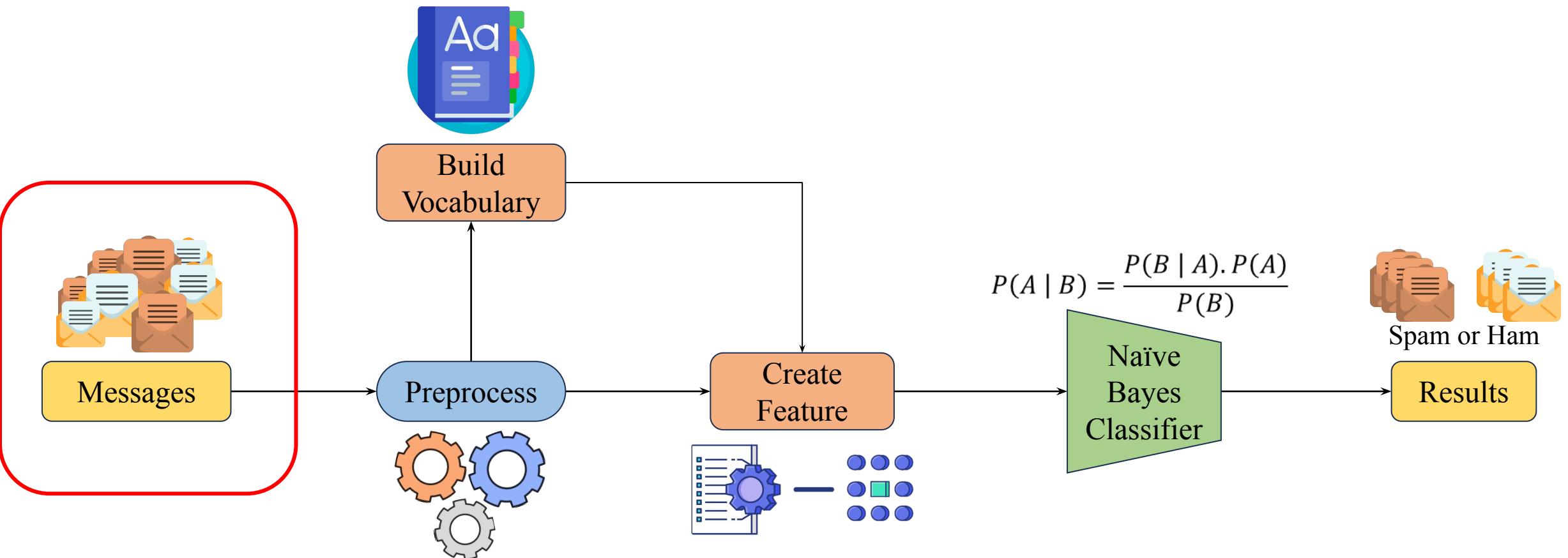
Code Implementation

◆ Project Pipeline



Code Implementation

◆ Project Pipeline



Code Implementation

❖ Coding step: Import libraries

```
1 import string
2 import nltk
3 nltk.download('stopwords')
4 nltk.download('punkt')
5 import pandas as pd
6 import numpy as np
7 import matplotlib.pyplot as plt
8
9 from sklearn.model_selection import train_test_split
10 from sklearn.naive_bayes import GaussianNB
11 from sklearn.metrics import accuracy_score
12 from sklearn.preprocessing import LabelEncoder
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]  Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]  Unzipping tokenizers/punkt.zip.
```



scikit-learn (sklearn): An open-source library for Python language that features various classification, regression and clustering algorithms.

Code Implementation

❖ Coding step: Download dataset

Description: Given Message Classification Dataset, build a Naives Bayes model to determine whether a text message is spam message or not (ham).

```
1 # https://drive.google.com/file/d/1N7rk-kfnDFIGMeX0R0VTjKh71gcgx-7R/view?usp=sharing
2 !gdown --id 1N7rk-kfnDFIGMeX0R0VTjKh71gcgx-7R
```

```
/usr/local/lib/python3.10/dist-packages/gdown/__main__.py:132: FutureWarning: Option `--id`
  warnings.warn(
Downloading...
From: https://drive.google.com/uc?id=1N7rk-kfnDFIGMeX0R0VTjKh71gcgx-7R
To: /content/2cls_spam_text_cls.csv
100% 486k/486k [00:00<00:00, 6.59MB/s]
```

Code Implementation

❖ Coding step: Read dataset

To read .csv file, we use pandas.read_csv():

```
1 DATASET_PATH = '/content/2cls_spam_text_cls.csv'  
2 df = pd.read_csv(DATASET_PATH)  
3 df
```

	Category	Message
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...
...
5567	spam	This is the 2nd time we have tried 2 contact u...
5568	ham	Will ü b going to esplanade fr home?
5569	ham	Pity, * was in mood for that. So...any other s...
5570	ham	The guy did some bitching but I acted like i'd...
5571	ham	Rofl. Its true to its name

5572 rows × 2 columns

Code Implementation

❖ Coding step: Read dataset

	Category	Message
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...
...
5567	spam	This is the 2nd time we have tried 2 contact u...
5568	ham	Will ü b going to esplanade fr home?
5569	ham	Pity, * was in mood for that. So...any other s...
5570	ham	The guy did some bitching but I acted like i'd...
5571	ham	Rofl. Its true to its name

5572 rows × 2 columns

```
1 messages = df['Message'].values.tolist()  
2 labels = df['Category'].values.tolist()
```

pandas.Series.to_list

`Series.to_list()`

[\[source\]](#)

Return a list of the values.

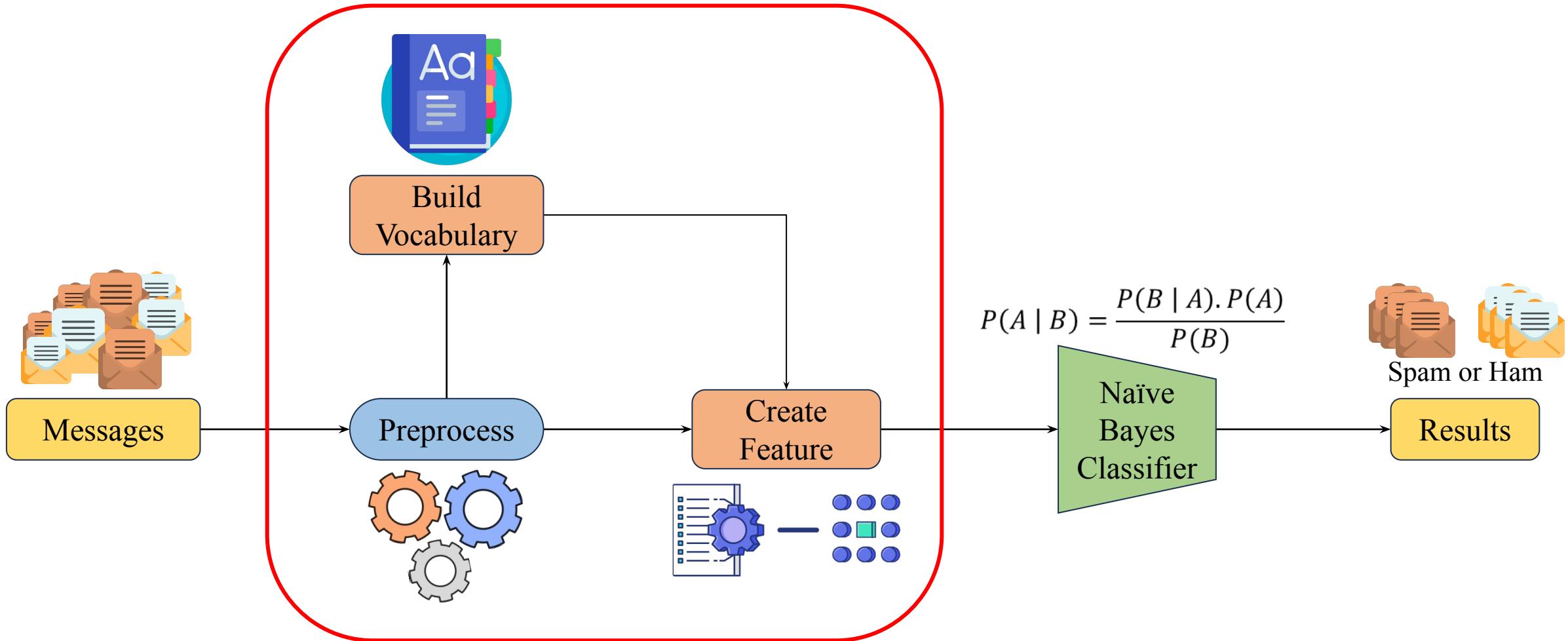
These are each a scalar type, which is a Python scalar (for str, int, float) or a pandas scalar (for Timestamp/Timedelta/Interval/Period)

Returns:

list

Code Implementation

❖ Project Pipeline



Code Implementation

◆ Text features

Corpus

doc1 = “deep learning book”

doc2 = “machine learning algorithm”

doc3 = “learning ai from scratch”

doc4 = “ai vietnam”

Tokenization



[‘deep’, ‘learning’, ‘book’]

[‘machine’, ‘learning’, ‘algorithm’]

[‘learning’, ‘ai’, ‘from’, ‘scratch’]

[‘ai’, ‘vietnam’]

Vocabulary

=

deep learning book machine algorithm ai from scratch vietnam

👉 Given a string = “vietnam machine learning deep learning book”

deep learning book machine algorithm ai from scratch vietnam

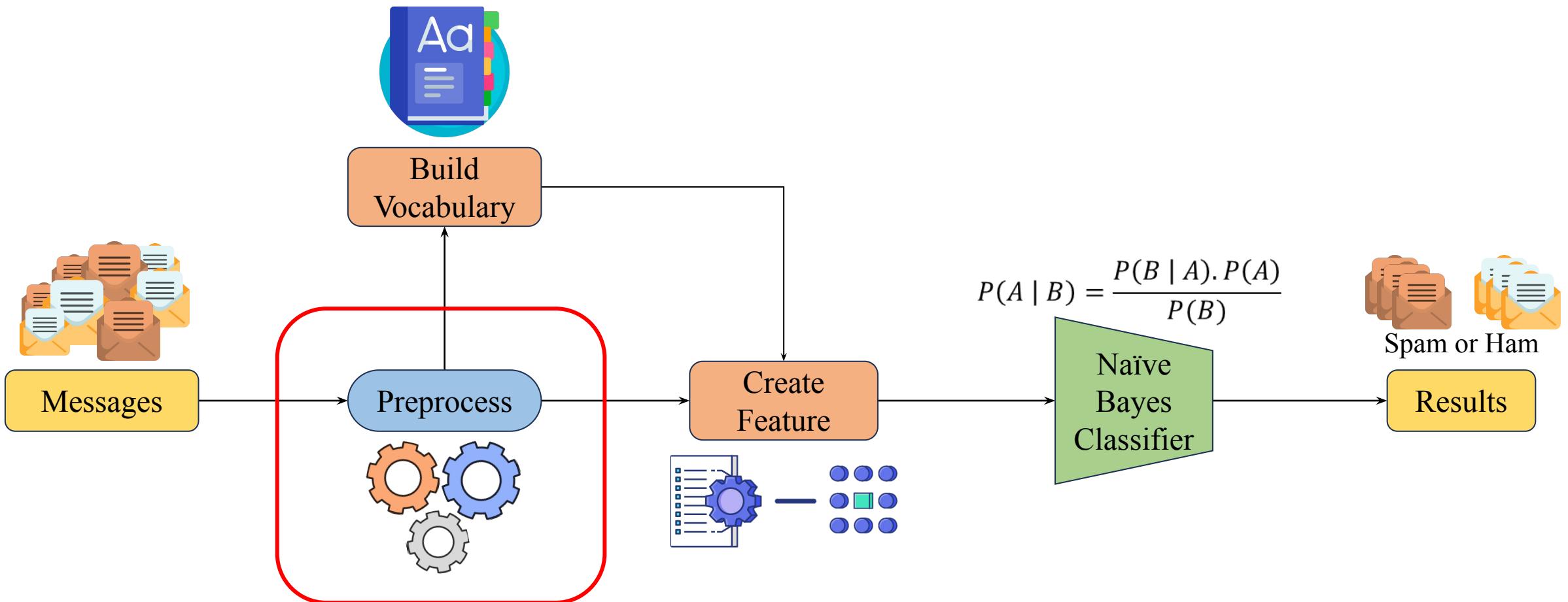
BoW 1 2 1 1 0 0 0 0 1

Binary BoW

1 1 1 1 0 0 0 0 1

Code Implementation

❖ Project Pipeline



Code Implementation

❖ Coding step: Data preprocessing

Category	Message
0 ham	Go until jurong point, crazy.. Available only ...
1 ham	Ok lar... Joking wif u oni...
2 spam	Free entry in 2 a wkly comp to win FA Cup fina...
3 ham	U dun say so early hor... U c already then say...
4 ham	Nah I don't think he goes to usf, he lives aro...
...	...
5567 spam	This is the 2nd time we have tried 2 contact u...
5568 ham	Will ü b going to esplanade fr home?
5569 ham	Pity, * was in mood for that. So...any other s...
5570 ham	The guy did some bitching but I acted like i'd...
5571 ham	Rofl. Its true to its name

5572 rows × 2 columns

Raw Text: We love this! Would you go? #talk
#makememories #unplug #relax #iphone #smartphone #wifi
#connect... http://fb.me/6N3LsUpCu

Data
Preprocessing

Vector Representation: array([4, 1, 5, 2, 0, 7, 9])

Code Implementation

◆ Data Preprocessing Mail Content

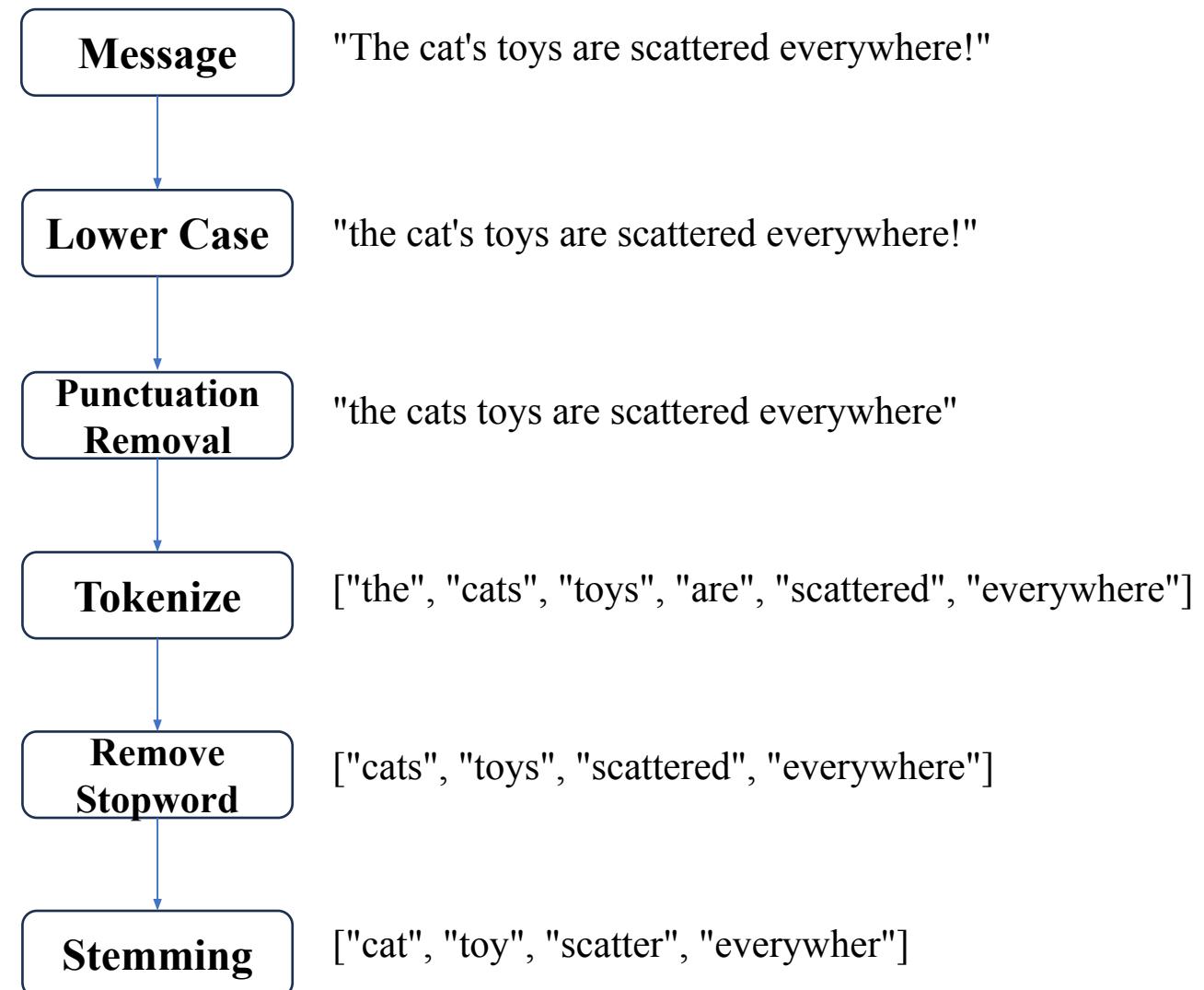
Converts all text to lowercase

Eliminates all punctuation marks

Splits the text into individual words (tokens)

Filters out common words that don't carry significant meaning

Reduces words to their root form, grouping similar words together

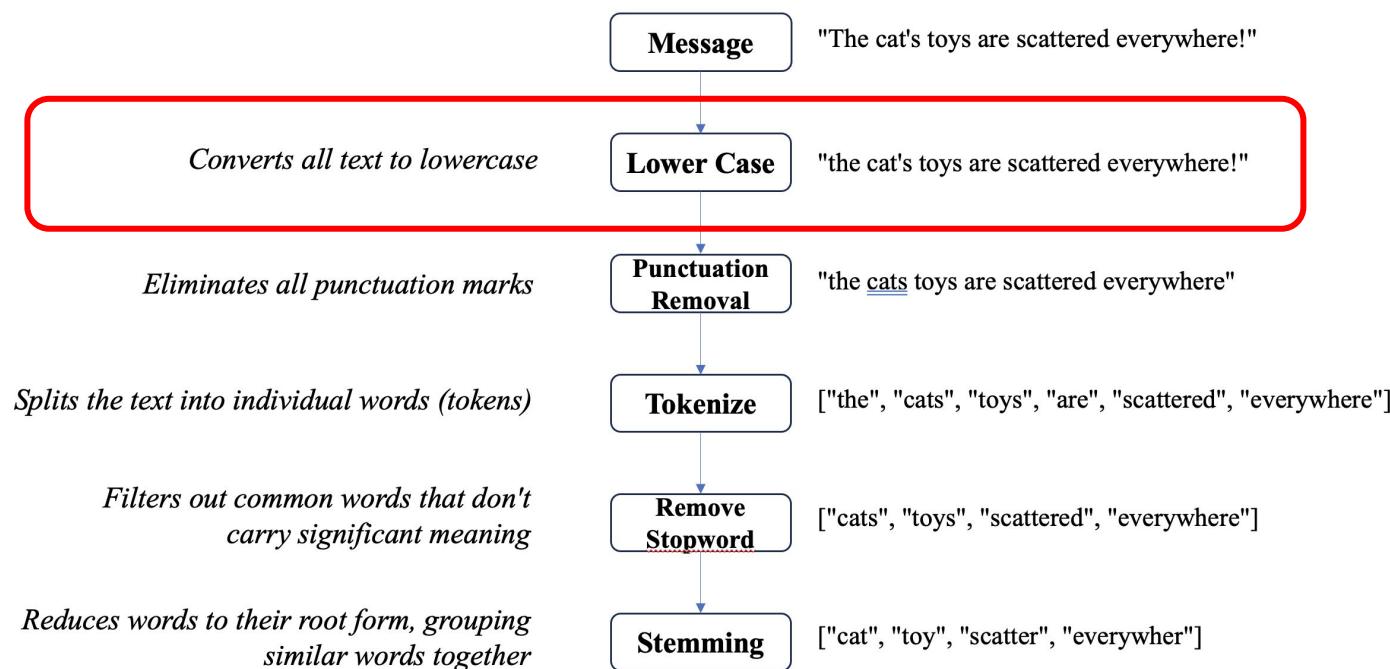


Code Implementation

❖ Data Preprocessing Mail Content

```
1 def lowercase(text):  
2     return text.lower()  
3  
4 INPUT_TEXT = "The cat's toys are scattered everywhere!"  
5 lowercase(INPUT_TEXT)
```

'the cat's toys are scattered everywhere!'

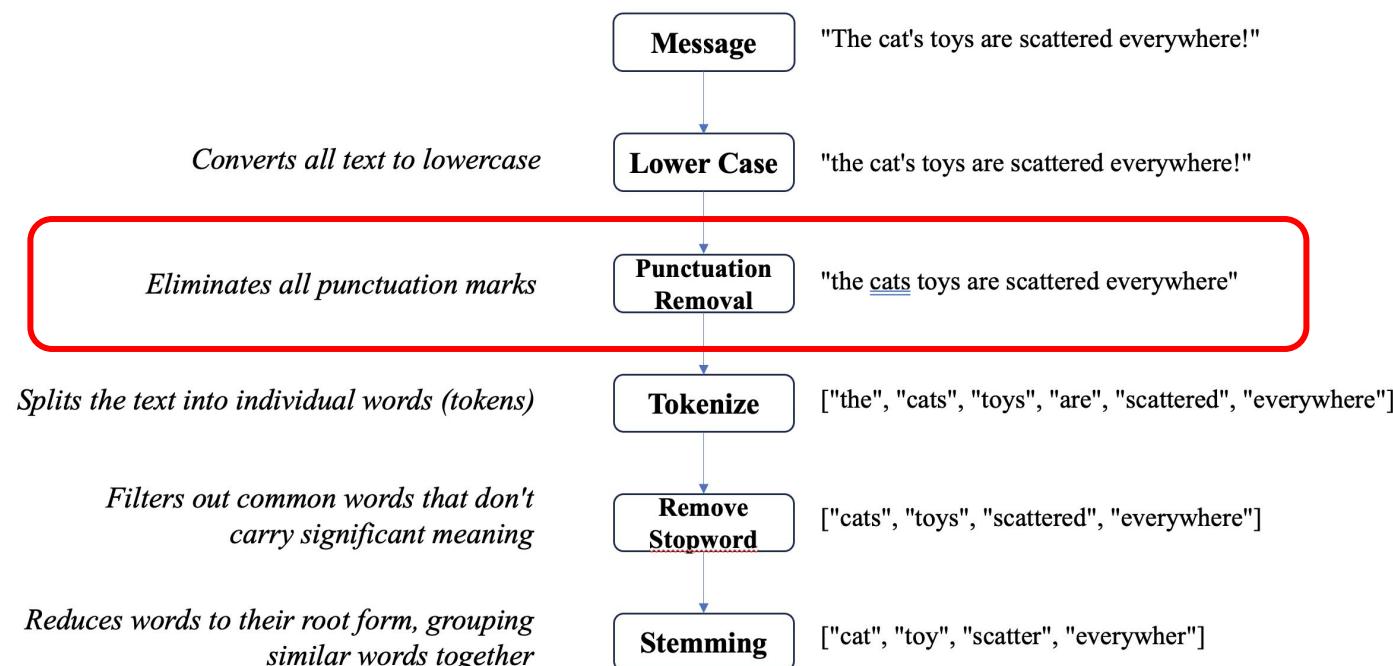


Code Implementation

◆ Data Preprocessing Mail Content

```
1 def punctuation_removal(text):  
2     translator = str.maketrans(' ', ' ', string.punctuation)  
3  
4     return text.translate(translator)  
5  
6 INPUT_TEXT = "The cat's toys are scattered everywhere!"  
7 INPUT_TEXT = lowercase(INPUT_TEXT)  
8 punctuation_removal(INPUT_TEXT)
```

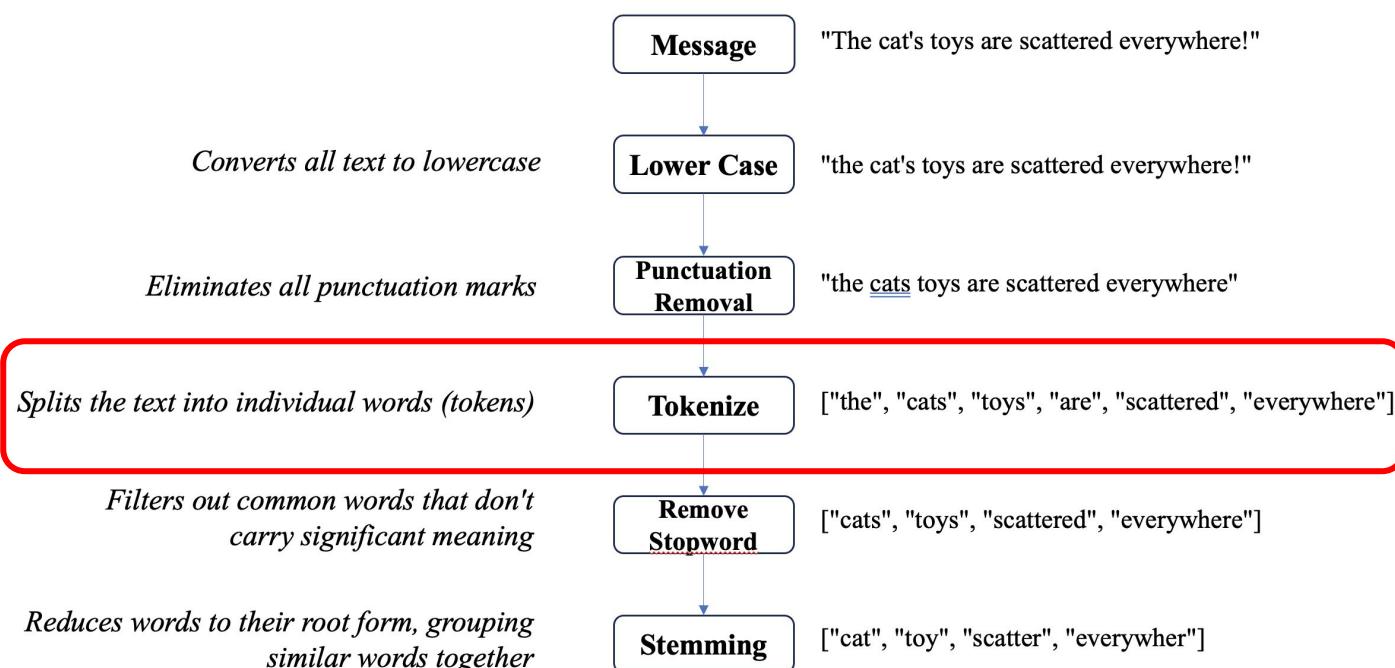
'the cats toys are scattered everywhere'



Code Implementation

◆ Data Preprocessing Mail Content

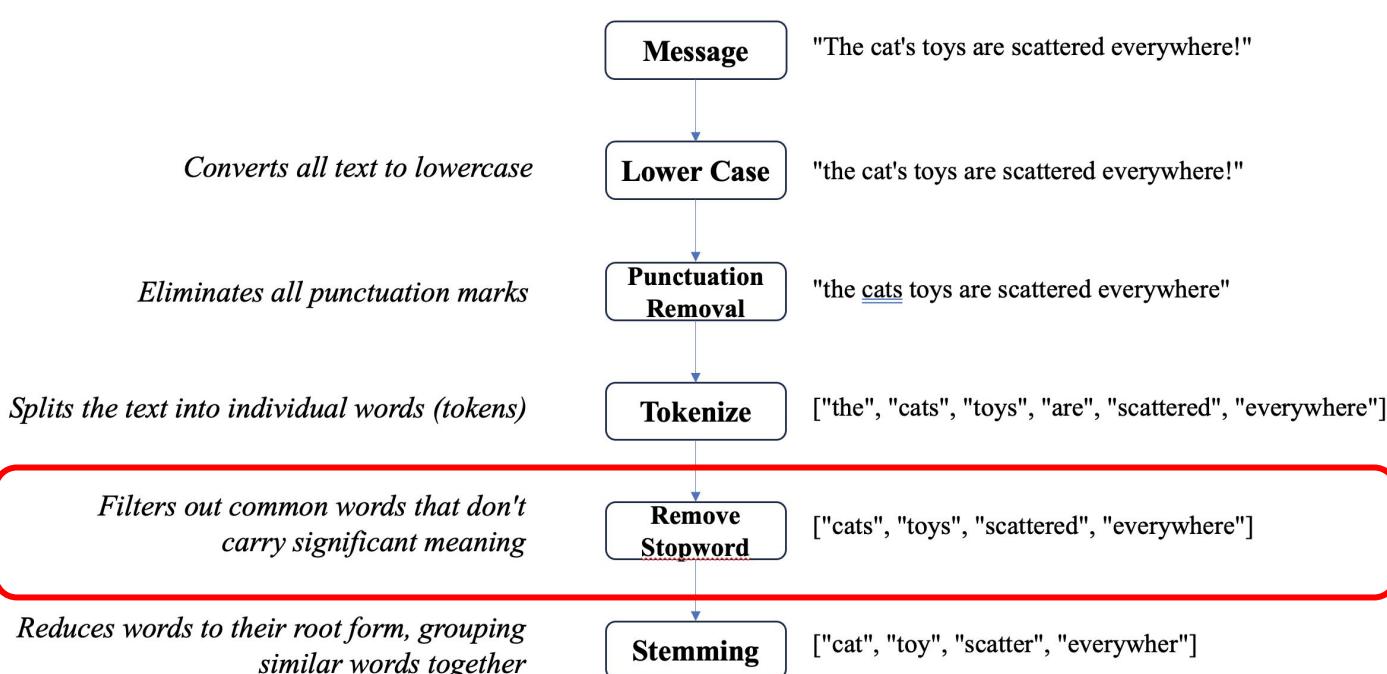
```
1 def tokenize(text):  
2     return nltk.word_tokenize(text)  
  
3  
4 INPUT_TEXT = "The cat's toys are scattered everywhere!"  
5 INPUT_TEXT = lowercase(INPUT_TEXT)  
6 INPUT_TEXT = punctuation_removal(INPUT_TEXT)  
7 tokenize(INPUT_TEXT)  
  
['the', 'cats', 'toys', 'are', 'scattered', 'everywhere']
```



Code Implementation

◆ Data Preprocessing Mail Content

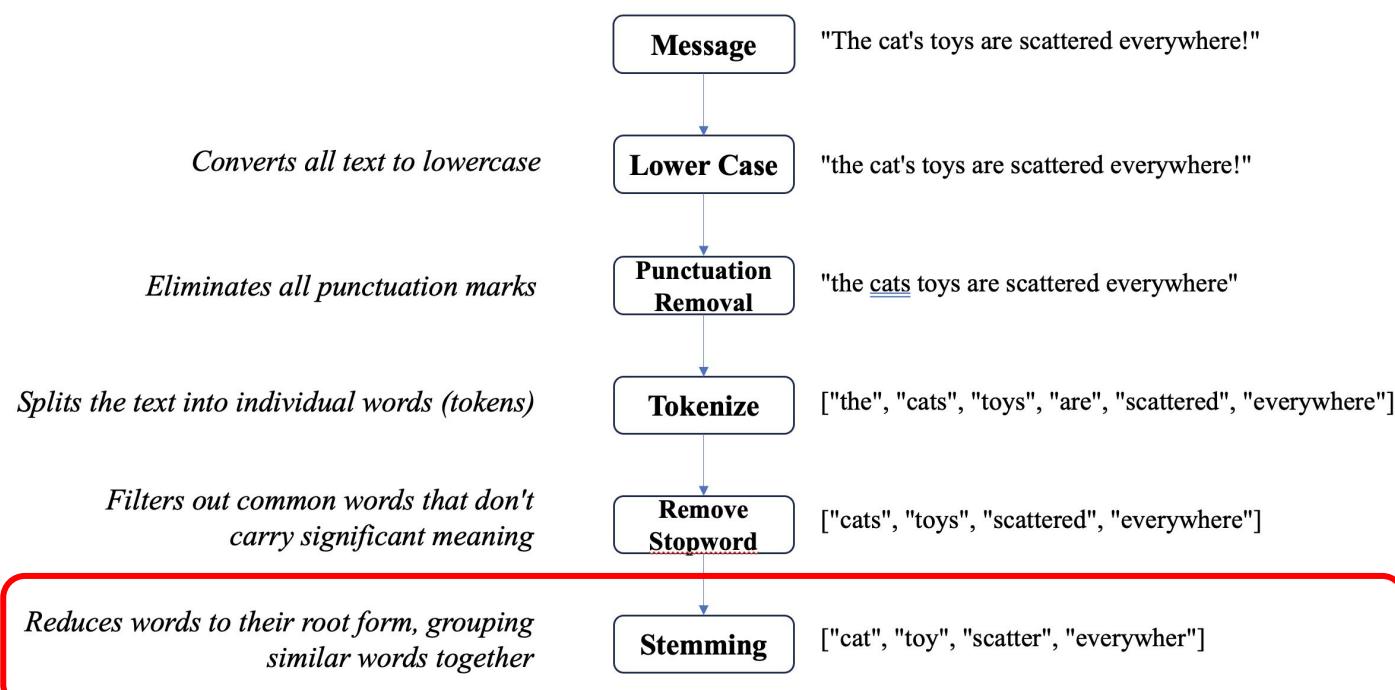
```
1 def remove_stopwords(tokens):  
2     stop_words = nltk.corpus.stopwords.words('english')  
3  
4     return [token for token in tokens if token not in stop_words]  
5  
6 INPUT_TEXT = "The cat's toys are scattered everywhere!"  
7 INPUT_TEXT = lowercase(INPUT_TEXT)  
8 INPUT_TEXT = punctuation_removal(INPUT_TEXT)  
9 INPUT_TEXT = tokenize(INPUT_TEXT)  
10 remove_stopwords(INPUT_TEXT)  
  
['cats', 'toys', 'scattered', 'everywhere']
```



Code Implementation

◆ Data Preprocessing Mail Content

```
1 def stemming(tokens):  
2     stemmer = nltk.PorterStemmer()  
3  
4     return [stemmer.stem(token) for token in tokens]  
5  
6 INPUT_TEXT = "The cat's toys are scattered everywhere!"  
7 INPUT_TEXT = lowercase(INPUT_TEXT)  
8 INPUT_TEXT = punctuation_removal(INPUT_TEXT)  
9 INPUT_TEXT = tokenize(INPUT_TEXT)  
10 INPUT_TEXT = remove_stopwords(INPUT_TEXT)  
11 stemming(INPUT_TEXT)  
  
['cat', 'toy', 'scatter', 'everywher']
```



Code Implementation

◆ Data Preprocessing Mail Content

```
1 def preprocess_text(text):
2     text = lowercase(text)
3     text = punctuation_removal(text)
4     tokens = tokenize(text)
5     tokens = remove_stopwords(tokens)
6     tokens = stemming(tokens)
7
8     return tokens
9
10 INPUT_TEXT = "The cat's toys are scattered everywhere!"
11 preprocess_text(INPUT_TEXT)
12
13 ['cat', 'toy', 'scatter', 'everywher']
```

Converts all text to lowercase

Eliminates all punctuation marks

Splits the text into individual words (tokens)

Filters out common words that don't carry significant meaning

Reduces words to their root form, grouping similar words together



Code Implementation

◆ Data Preprocessing Mail Content

Category	Message
0 ham	Go until jurong point, crazy.. Available only ...
1 ham	Ok lar... Joking wif u oni...
2 spam	Free entry in 2 a wkly comp to win FA Cup fina...
3 ham	U dun say so early hor... U c already then say...
4 ham	Nah I don't think he goes to usf, he lives aro...
...	...
5567 spam	This is the 2nd time we have tried 2 contact u...
5568 ham	Will ü b going to esplanade fr home?
5569 ham	Pity, * was in mood for that. So...any other s...
5570 ham	The guy did some bitching but I acted like i'd...
5571 ham	Rofl. Its true to its name

5572 rows × 2 columns

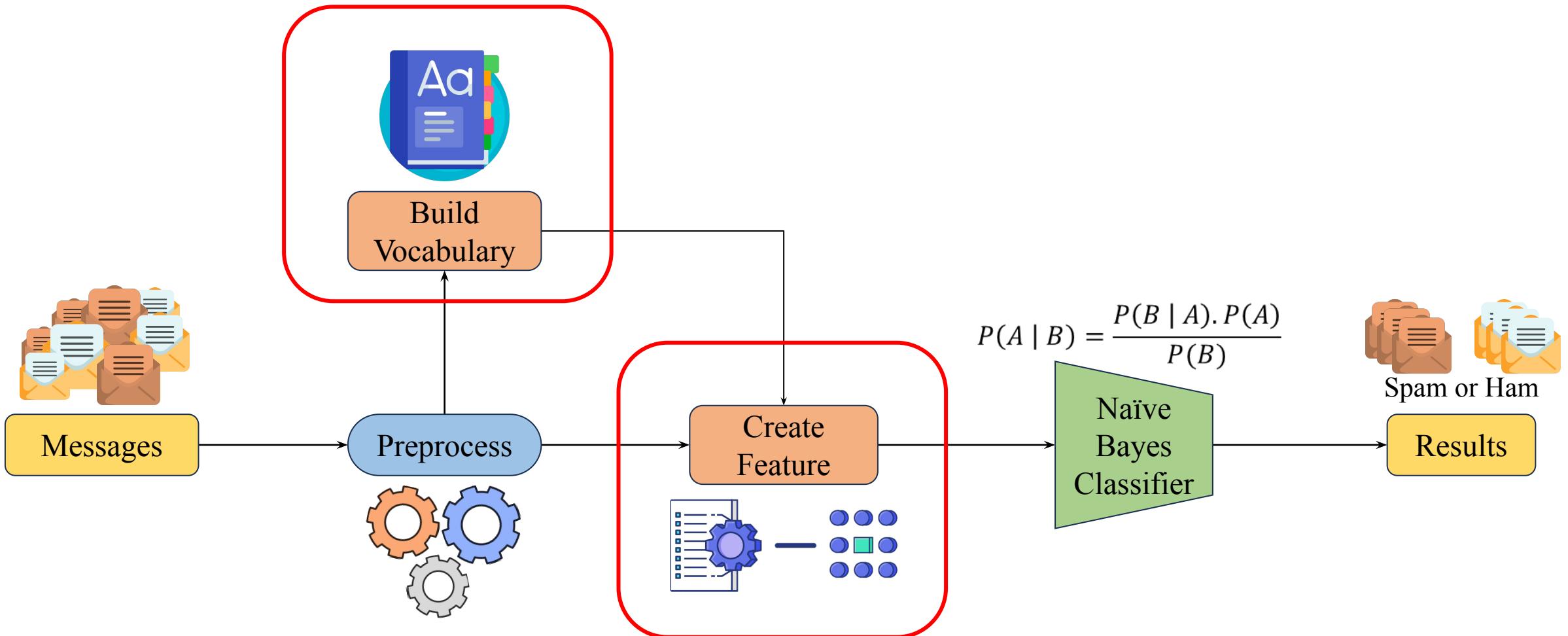
```
1 messages = [preprocess_text(message) for message in messages]
```

1 messages

```
[['go',  
 'jurong',  
 'point',  
 'crazi',  
 'avail',  
 'bugi',  
 'n',  
 'great',  
 'world',  
 'la',  
 'e',  
 'buffet',  
 'cine',  
 'got',  
 'amor',  
 'wat'],  
 ['ok', 'lar', 'joke', 'wif', 'u', 'oni'],
```

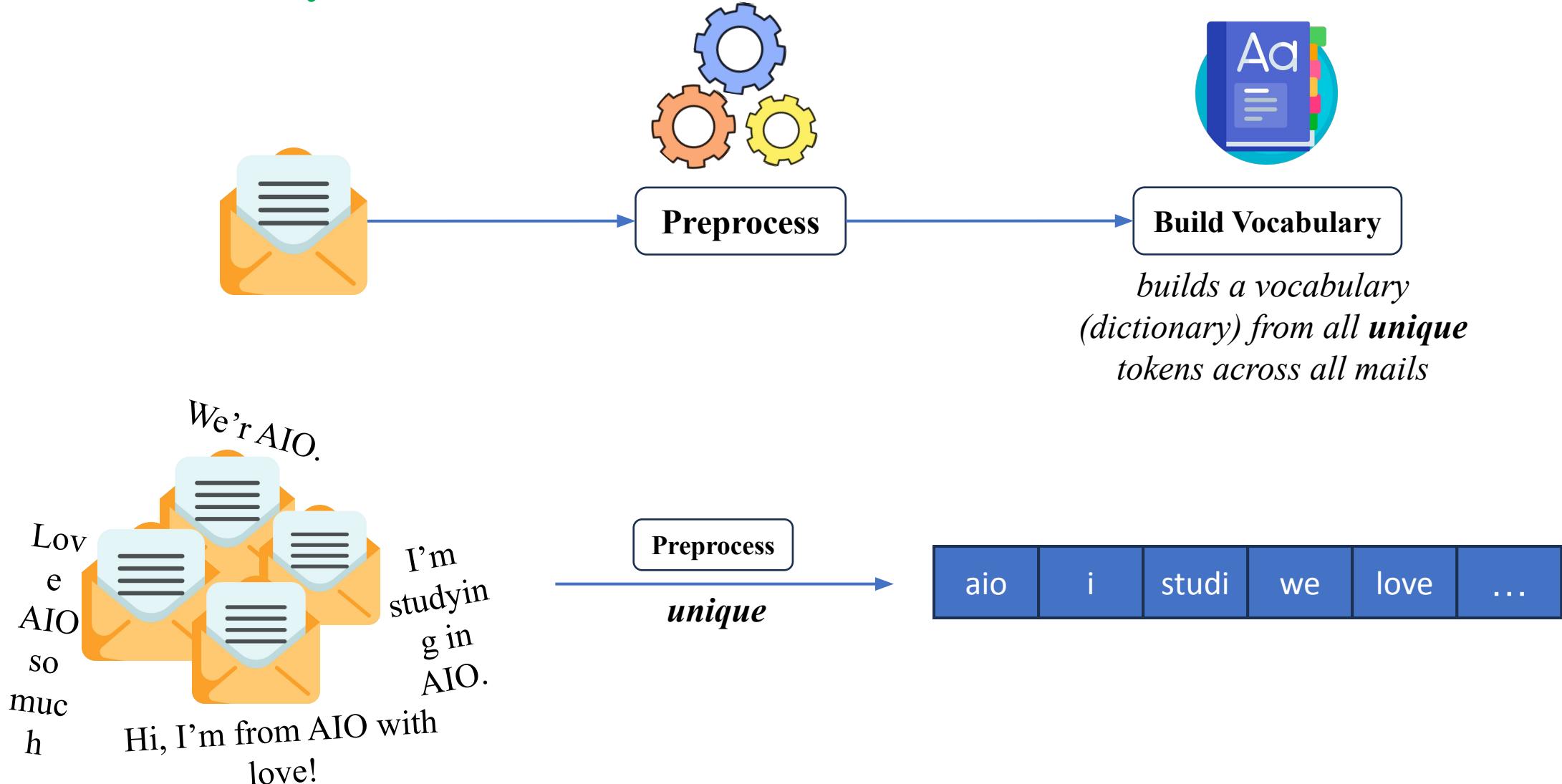
Code Implementation

◆ Project Pipeline



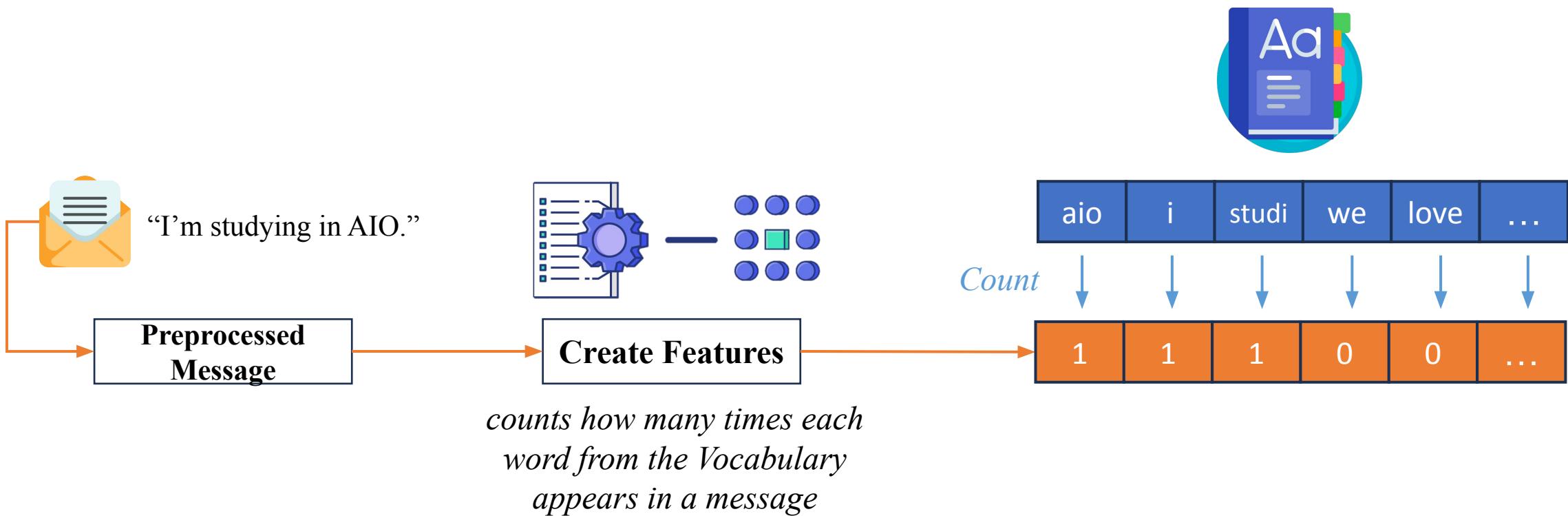
Code Implementation

❖ Build Dictionary



Code Implementation

❖ Create Features



Code Implementation

❖ Coding step: Build vocab

```
1 def create_dictionary(messages):
2     dictionary = []
3
4     for tokens in messages:
5         for token in tokens:
6             if token not in dictionary:
7                 dictionary.append(token)
8
9     return dictionary
10
11 dictionary = create_dictionary(messages)
12 dictionary
```

```
['go',
'jurong',
'point',
'crazi',
'avail',
'bugi',
'n',
'great',
'world',
```



Build Vocabulary

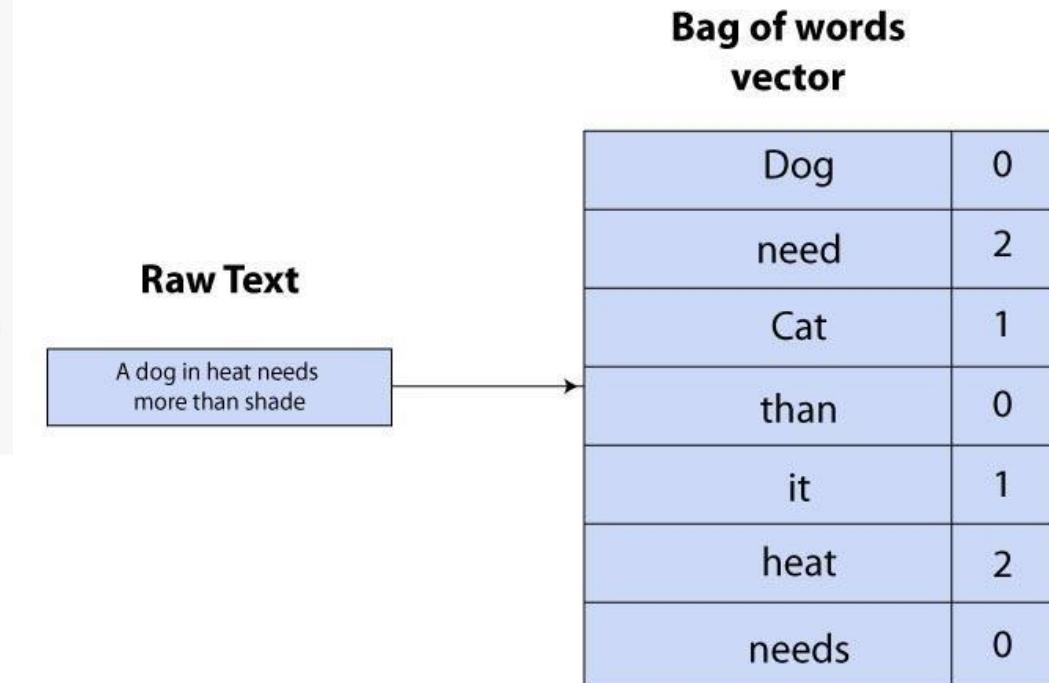
*builds a vocabulary
(dictionary) from all **unique**
tokens across all mails*

Code Implementation

❖ Coding step: Create features

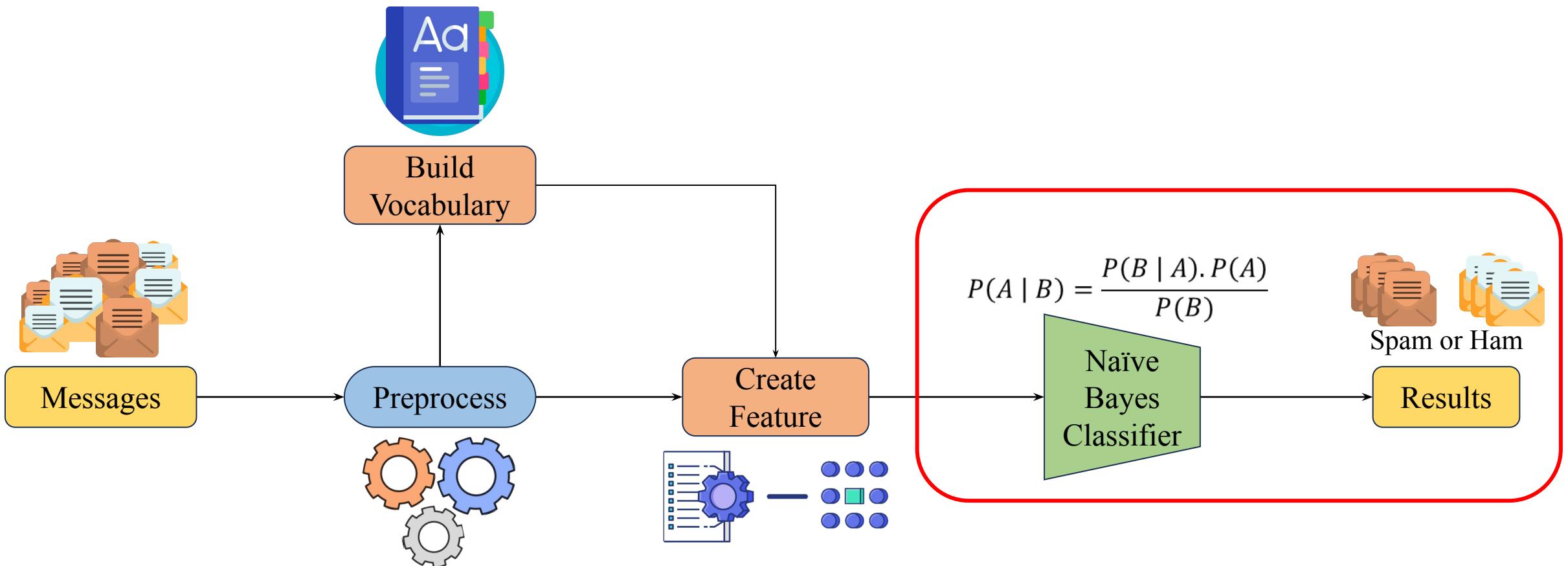
```
1 def create_features(tokens, dictionary):
2     features = np.zeros(len(dictionary))
3
4     for token in tokens:
5         if token in dictionary:
6             features[dictionary.index(token)] += 1
7
8     return features
9
10 dictionary = create_dictionary(messages)
11 X = np.array([create_features(tokens, dictionary) for tokens in messages])
12 print(X.shape)
13 print(X)

(5572, 8166)
[[1. 1. 1. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 ...
 [0. 0. 0. ... 0. 1. 1.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]]
```



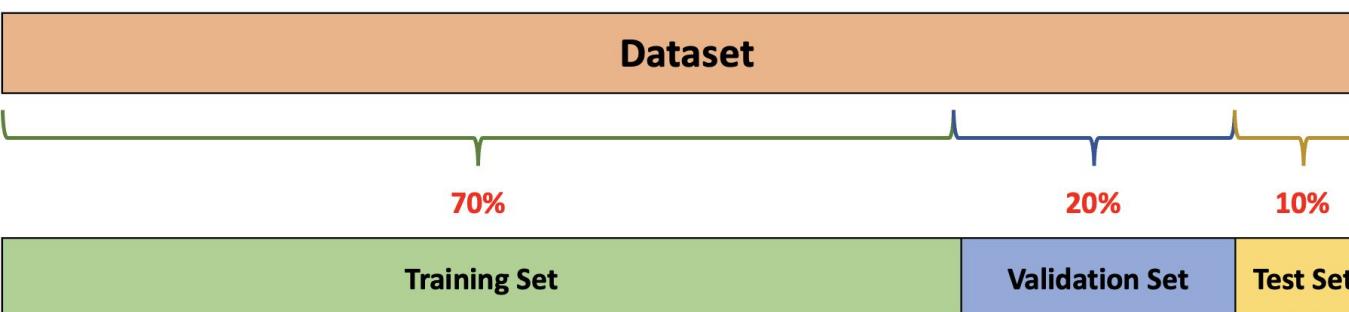
Code Implementation

◆ Project Pipeline



Code Implementation

❖ Coding step: Split train, val, test set



```
1 VAL_SIZE = 0.2
2 TEST_SIZE = 0.125
3 SEED = 0
4 IS_SHUFFLE = True
5
6 X_train, X_val, y_train, y_val = train_test_split(
7     X, y,
8     test_size=VAL_SIZE,
9     shuffle=IS_SHUFFLE,
10    random_state=SEED
11 )
12 X_train, X_test, y_train, y_test = train_test_split(
13     X_train, y_train,
14     test_size=TEST_SIZE,
15     shuffle=IS_SHUFFLE,
16     random_state=SEED
17 )
```

```
1 print(f'Number of training samples: {X_train.shape[0]}')
2 print(f'Number of val samples: {X_val.shape[0]}')
3 print(f'Number of test samples: {X_test.shape[0]}')
```

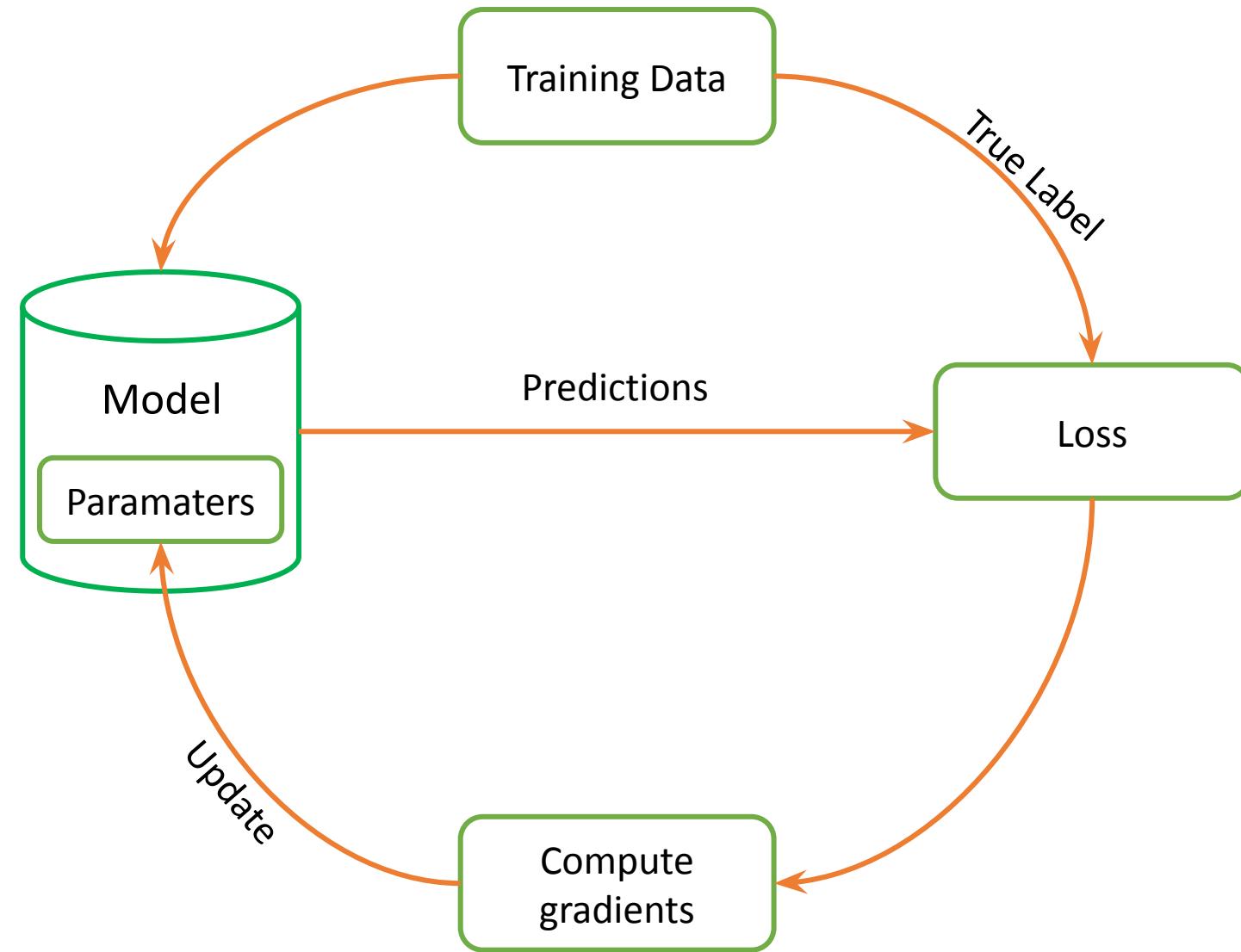
Number of training samples: 3899
Number of val samples: 1115
Number of test samples: 558

Code Implementation

❖ Coding step: Training model

```
1 %%time
2 model = GaussianNB()
3 print('Start training...')
4 model = model.fit(X_train, y_train)
5 print('Training completed!')
```

Start training...
Training completed!
CPU times: user 397 ms, sys: 162 ms, total: 559 ms
Wall time: 633 ms



Code Implementation

❖ Coding step: Evaluation

$$\text{accuracy} = \frac{\text{true_predictions}}{\text{n_samples}}$$

```
1 y_val_pred = model.predict(X_val)
2 y_test_pred = model.predict(X_test)
3 val_accuracy = accuracy_score(y_val, y_val_pred)
4 test_accuracy = accuracy_score(y_test, y_test_pred)
5 print(f'Val accuracy: {val_accuracy}')
6 print(f'Test accuracy: {test_accuracy}')
```

Val accuracy: 0.8816143497757848
Test accuracy: 0.8602150537634409

Code Implementation

❖ Coding step: Prediction

```
1 def predict(text, model, dictionary):
2     processed_text = preprocess_text(text)
3     features = create_features(text, dictionary)
4     features = np.array(features).reshape(1, -1)
5     prediction = model.predict(features)
6     prediction_cls = le.inverse_transform(prediction)[0]
7
8     return prediction_cls
```

```
1 test_input = 'I am actually thinking a way of doing something useful'
2 prediction_cls = predict(test_input, model, dictionary)
3 print(f'Prediction: {prediction_cls}')
```

Prediction: ham

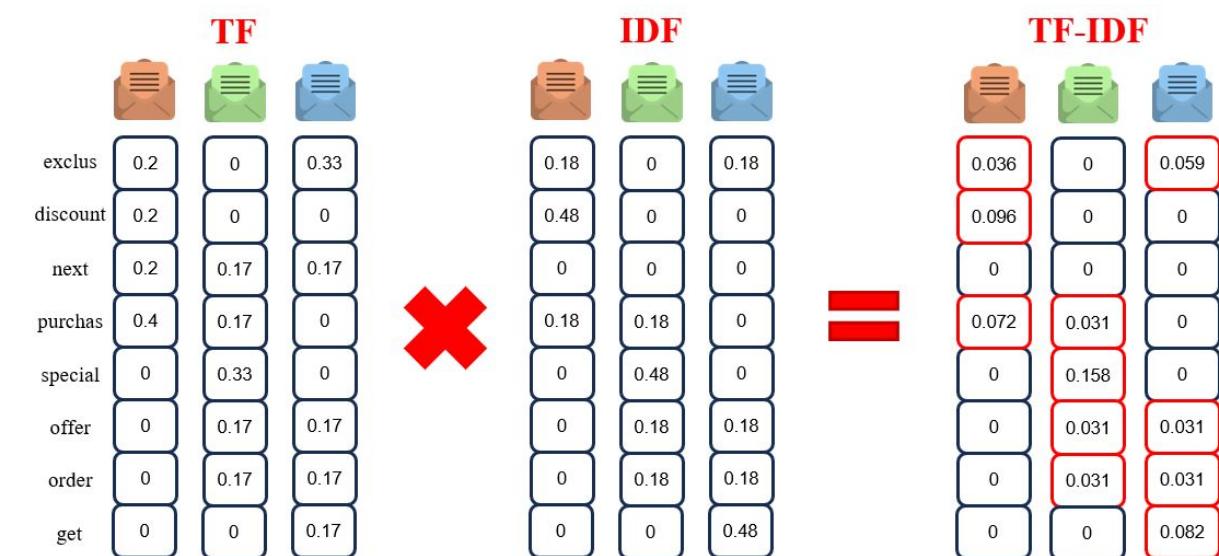
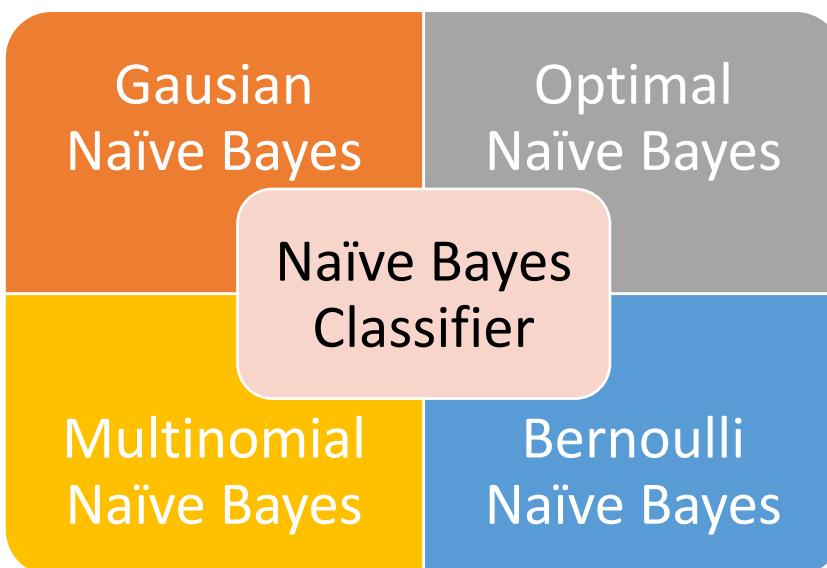
QUIZ

Improved Baseline

Improved Baseline

◆ Introduction

1. Use other Naive Bayes variables.
2. Use advanced feature extraction techniques such as TF-IDF, BoW, N-grams, ...



Improved Baseline

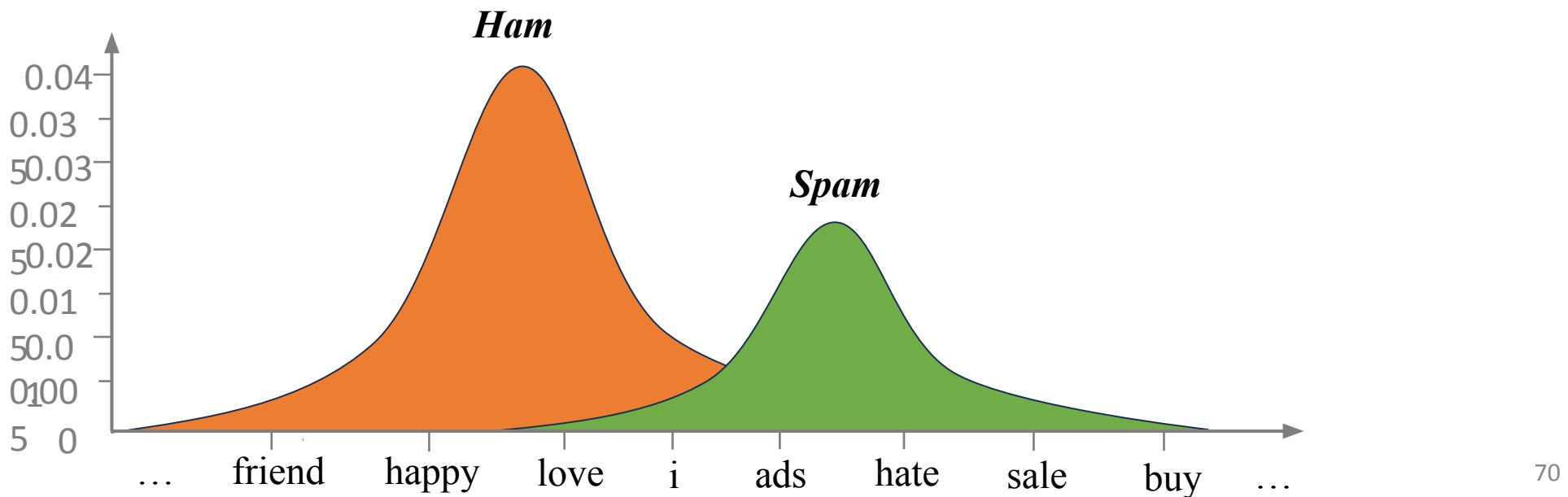
◆ 1. Different Naïve Bayes Classifier

The different naive Bayes classifiers differ mainly by the assumptions $P(x_i | y)$.

Gaussian Naive Bayes assumes that:

*Gaussian
Naïve Bayes
Model*

$$P(x_i | y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right), \text{ where mean } \mu \text{ and var } \sigma \text{ are estimated.}$$

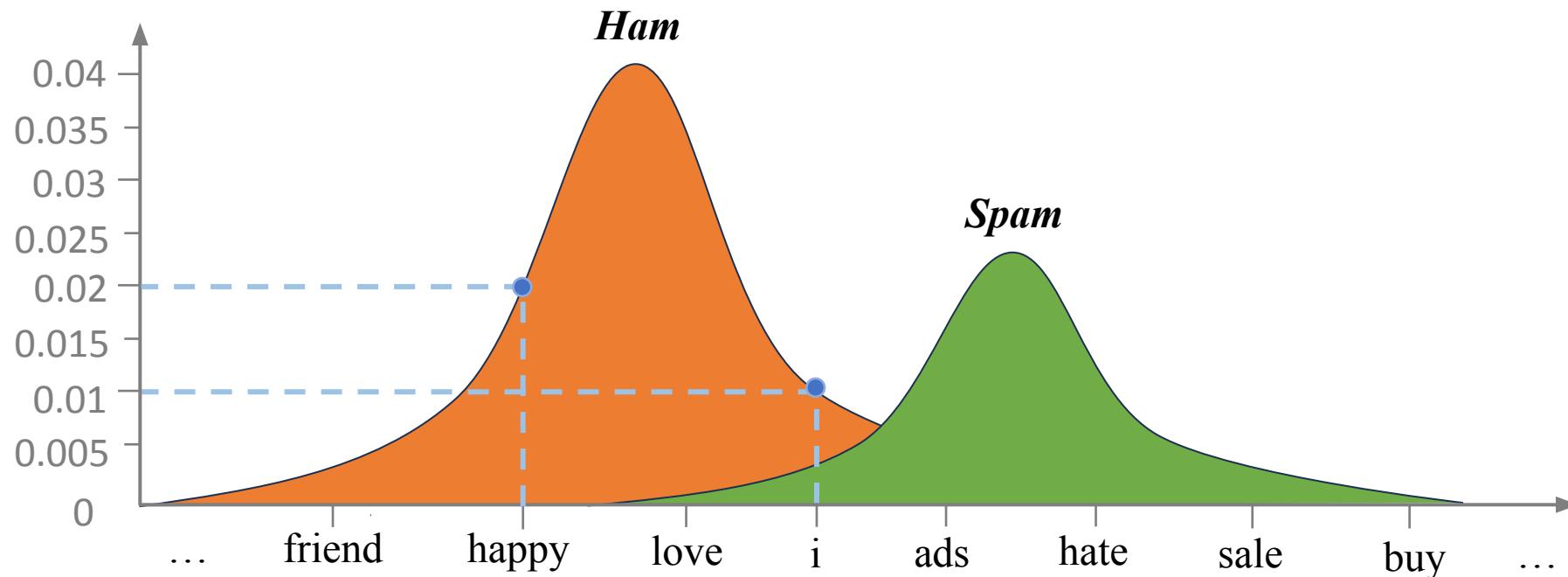


Improved Baseline

◆ 1. Different Naïve Bayes Classifier

New Message A: “*I am happy*”

$$P(A | Ham) = P("i" | Ham) \cdot P("happy" | Ham) = 0.01 * 0.02 = 0.002$$

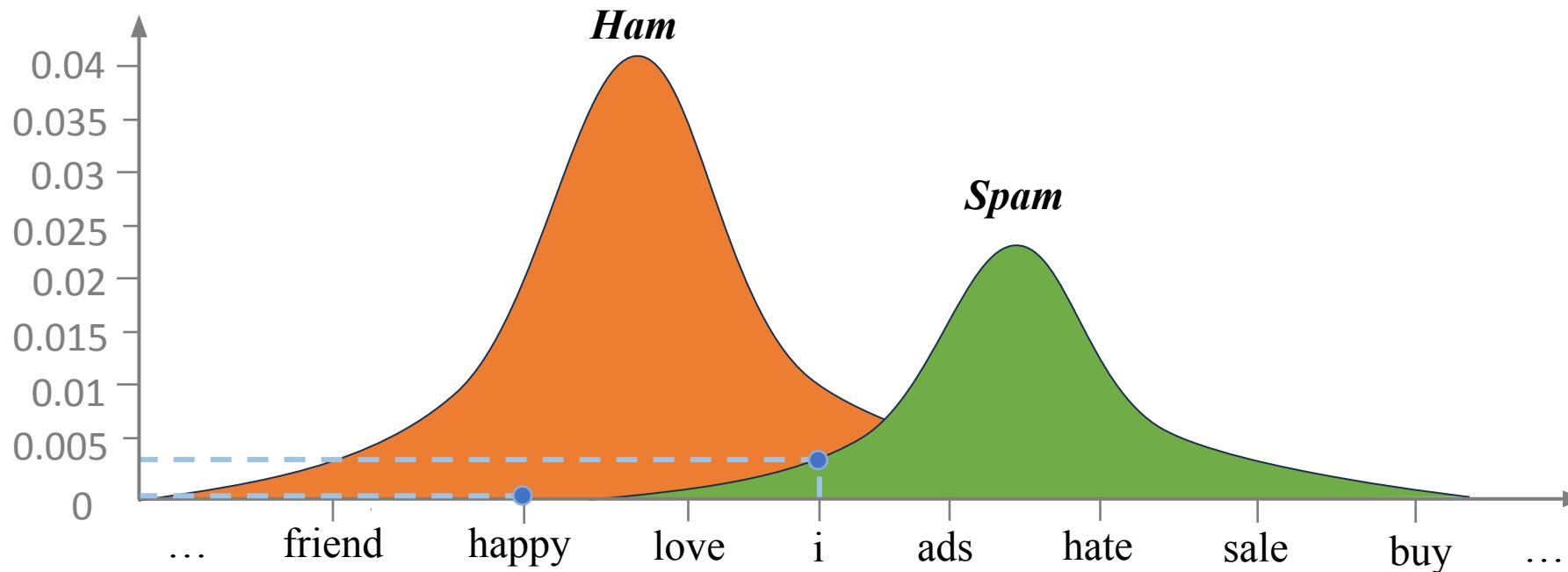


Improved Baseline

◆ 1. Different Naïve Bayes Classifier

New Message A: “*I am happy*”

$$P(A | \text{Spam}) = P("i"|\text{Spam}) \cdot P("happy"|\text{Spam}) = 0.003 * 0.00001 = 0.0000003$$

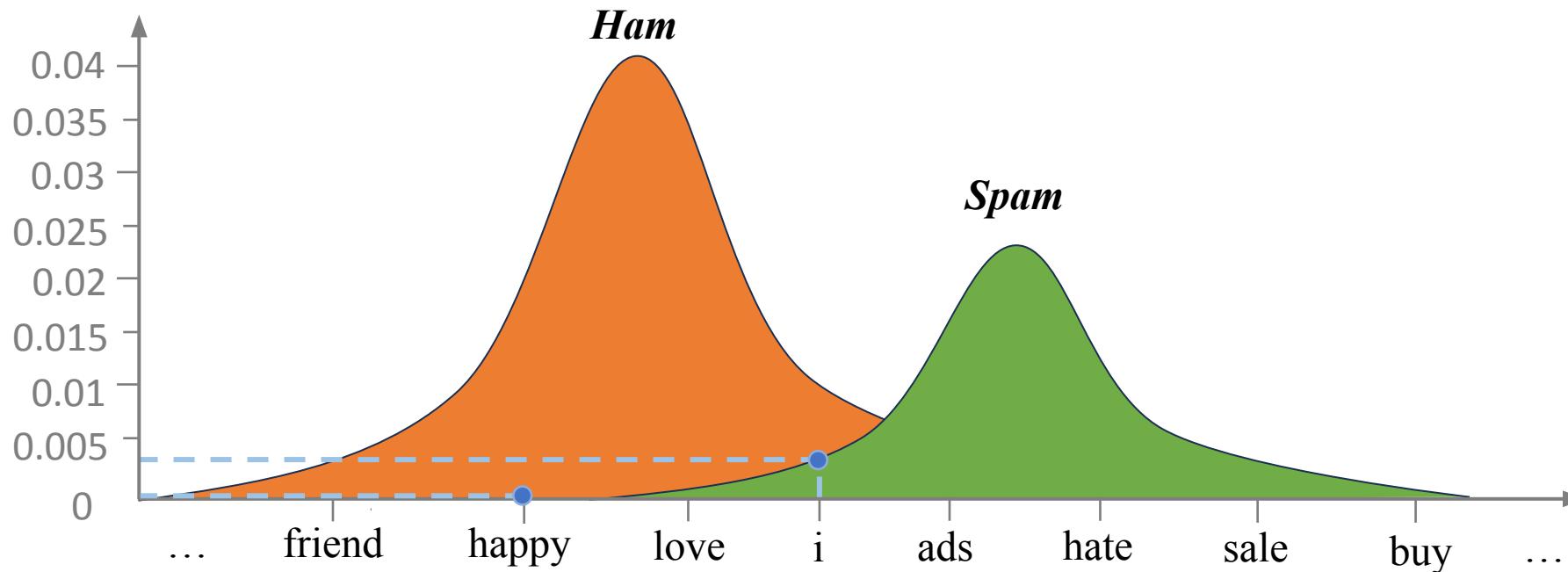


Improved Baseline

◆ 1. Different Naïve Bayes Classifier

New Message A: “*I am happy*”

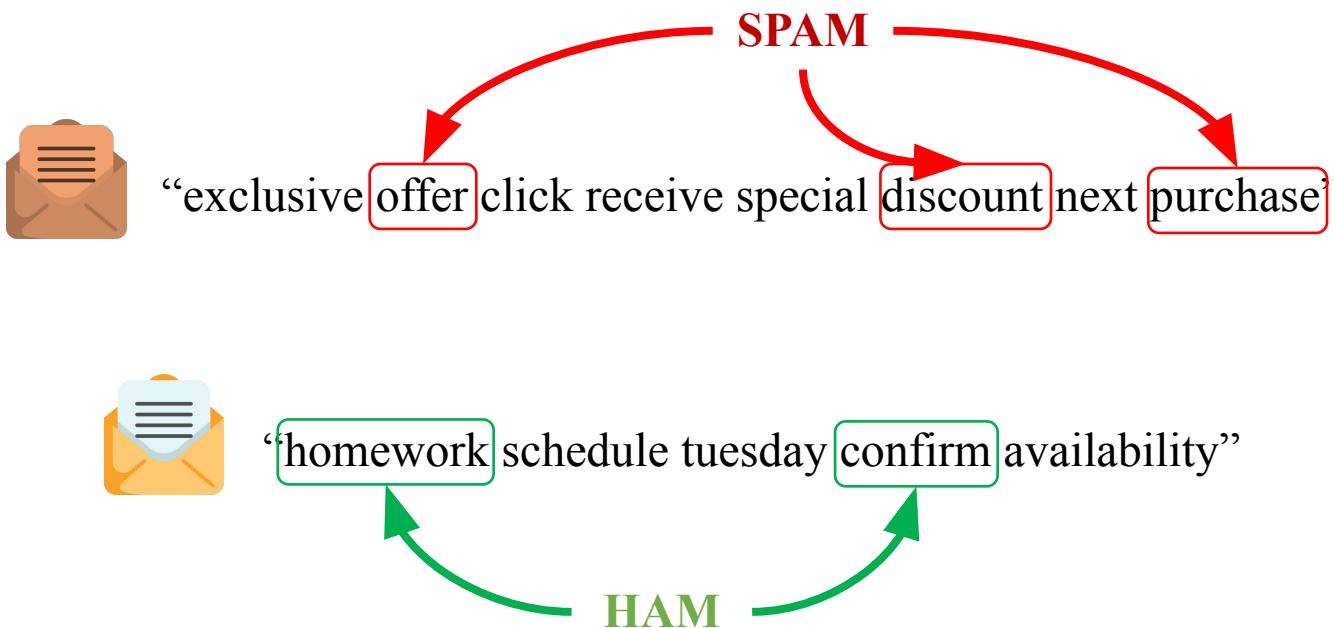
$$P(A | \text{Spam}) = P("i"|\text{Spam}) \cdot P("happy"|\text{Spam}) = 0.003 * 0.00001 = 0.0000003$$



Improved Baseline

◆ 2. Different Feature Extraction Techniques

TF-IDF (Term Frequency-Inverse Document Frequency), measure used to evaluate the importance of a word in a document relative to a collection of documents (corpus), to identify words that are significant in a specific document but not too common across all documents.



Improved Baseline

◆ 2. Different Feature Extraction Techniques

Term Frequency measures how frequently a word appears in a document.

$$TF = \frac{\text{Number of times } \mathbf{word}_i \text{ appears in mail}}{\text{Total number of words in mail}}$$



“exclus discount next purchas purchas”



“special offer next purchas special order”



“get exclus offer next exclus order”



$$\text{exclus} \quad 0.2 = \frac{1}{5}$$



$$\text{discount} \quad 0.2 \quad 0 \quad 0.3 \\ 3$$



$$\text{next} \quad 0.2 \quad 0 \quad 0 \\ 0.1 \quad 7 \quad 0.1 \\ 7 \quad 7$$

$$\text{purchas} \quad 0.4 = \frac{2}{5} \quad 0.1 \quad 0 \\ 7 \quad 0.1 \\ 7 \quad 0$$

$$\text{special} \quad 0 \quad 0.3 \\ 3 \quad 0$$

$$\text{offer} \quad 0 \quad 0.1 \\ 7 \quad 0.1 \\ 7 \quad 7$$

$$\text{order} \quad 0 \quad 0.1 \\ 7 \quad 0.1 \\ 7 \quad 7$$

$$\text{get} \quad 0 \quad 0.1 \\ 7 \quad 7$$

Improved Baseline

◆ 2. Different Feature Extraction Techniques

Inverse Document Frequency measures how important a word is across all mail.

$$IDF = \log \frac{\text{Number of mails}}{\text{Number of mails containing word}_i}$$



“exclus discount next purchas purchas”



“special offer next purchas special order”



“get exclus offer next exclus order”

exclus	<table border="1"><tr><td>0.1</td><td></td></tr><tr><td>8</td><td></td></tr></table>	0.1		8		$= \log \frac{3}{2}$	
0.1							
8							
discount	<table border="1"><tr><td>0.4</td><td></td></tr><tr><td>8</td><td></td></tr></table>	0.4		8		$= \log \frac{3}{1}$	
0.4							
8							
next	<table border="1"><tr><td>0</td><td></td></tr><tr><td></td><td></td></tr></table>	0				$= \log \frac{3}{3}$	
0							
purchas	<table border="1"><tr><td>0.1</td><td></td></tr><tr><td>8</td><td></td></tr></table>	0.1		8			
0.1							
8							
special	<table border="1"><tr><td>0</td><td></td></tr><tr><td></td><td></td></tr></table>	0					
0							
offer	<table border="1"><tr><td>0</td><td></td></tr><tr><td></td><td></td></tr></table>	0					
0							
order	<table border="1"><tr><td>0</td><td></td></tr><tr><td></td><td></td></tr></table>	0					
0							
get	<table border="1"><tr><td>0</td><td></td></tr><tr><td></td><td></td></tr></table>	0					
0							

Improved Baseline

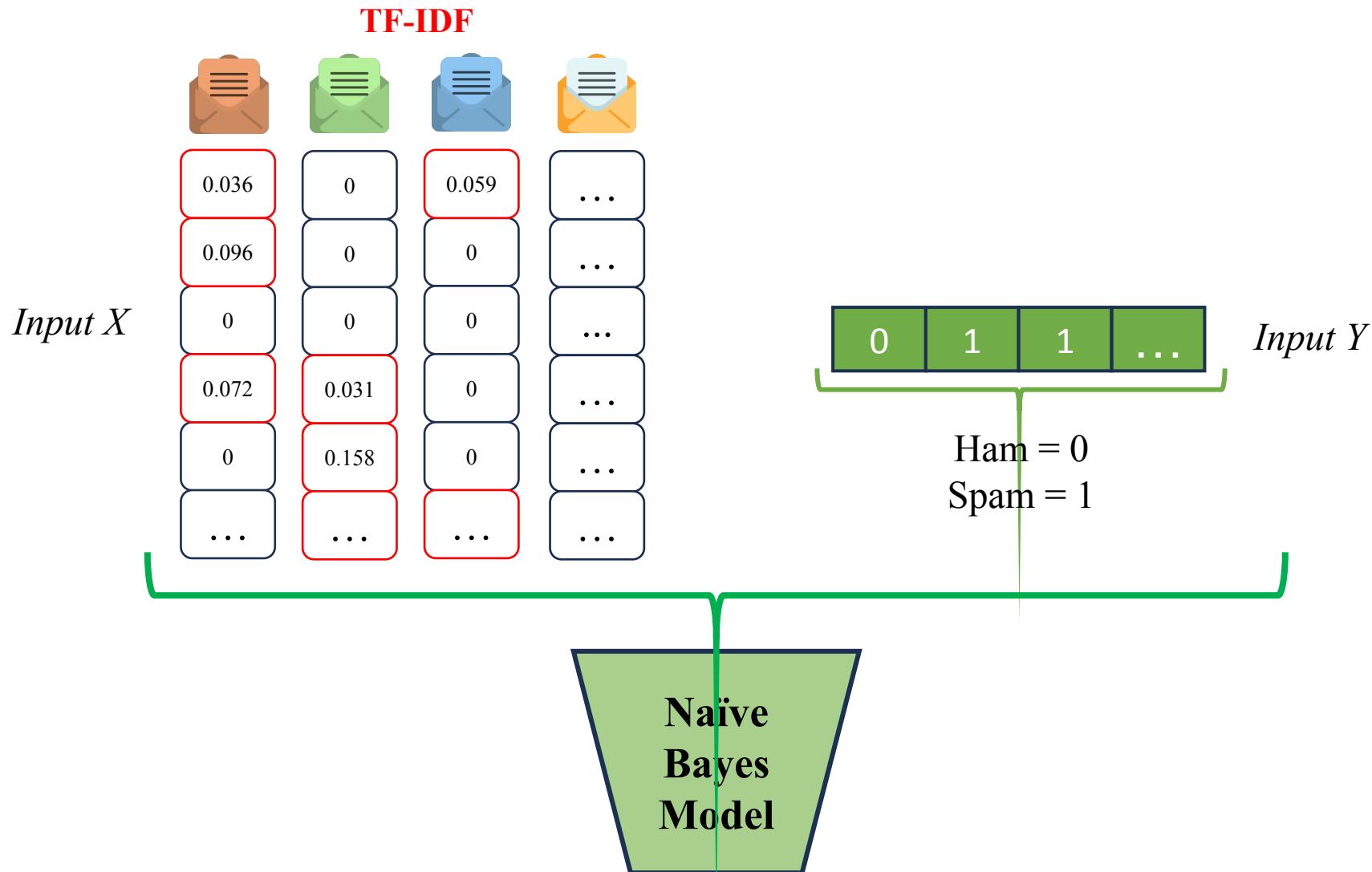
◆ 2. Different Feature Extraction Techniques

Combines TF and IDF to determine the importance of a word in a mail relative to the entire mails.

	TF			IDF			TF-IDF		
exclus	0.2	0	0.33	0.18	0	0.18	0.036	0	0.059
discount	0.2	0	0	0.48	0	0	0.096	0	0
next	0.2	0.17	0.17	0	0	0	0	0	0
purchas	0.4	0.17	0	0.18	0.18	0	0.072	0.031	0
special	0	0.33	0	0	0.48	0	0	0.158	0
offer	0	0.17	0.17	0	0.18	0.18	0	0.031	0.031
order	0	0.17	0.17	0	0.18	0.18	0	0.031	0.031
get	0	0	0.17	0	0	0.48	0	0	0.082

Improved Baseline

◆ Input Naïve Bayes Model



Summarization

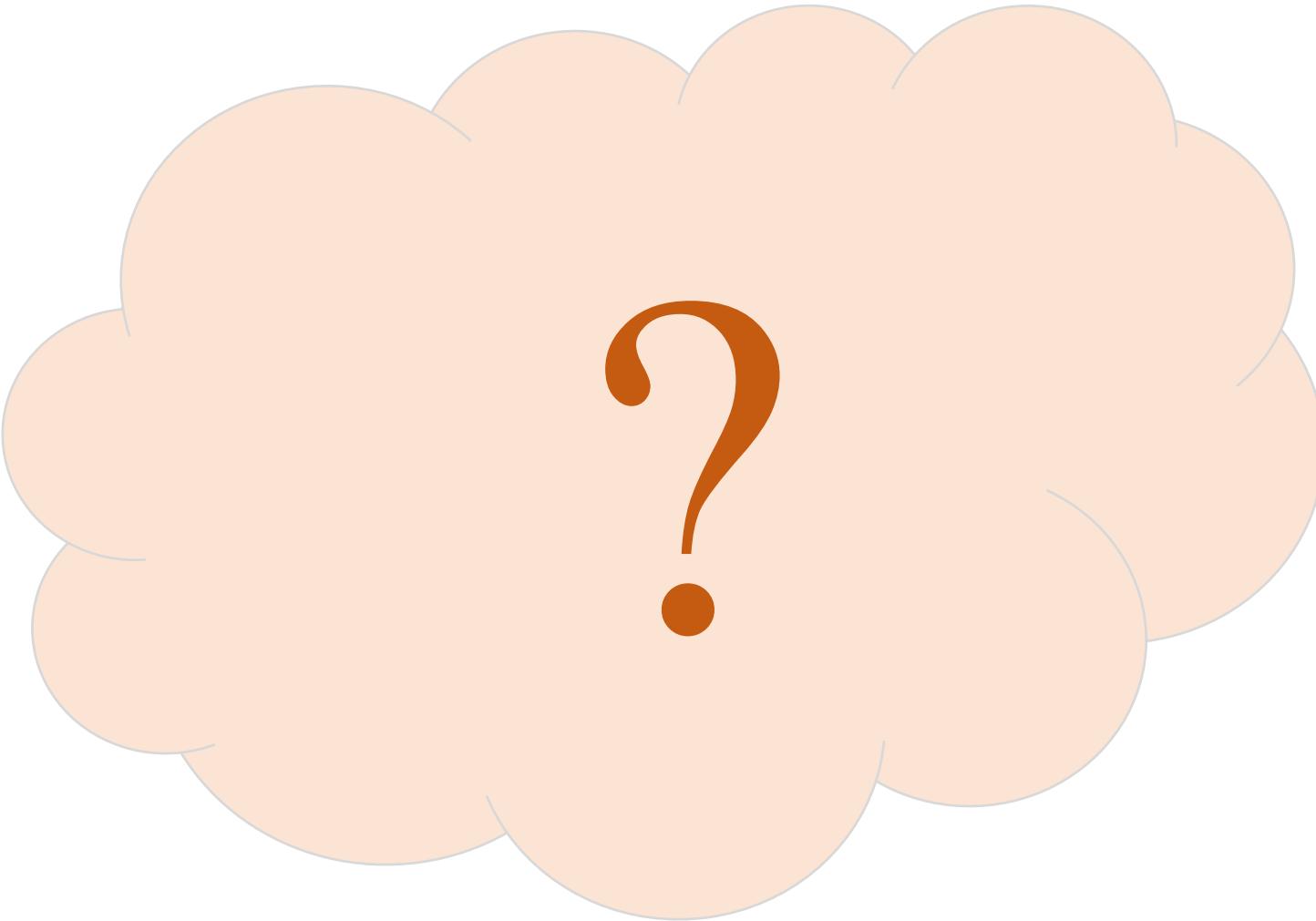
❖ What we have learned so far

Spam	Not spam
 <p>"SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply Reply HL 4 info"</p>	 <p>"Nah I don't think he goes to usf, he lives around here though"</p>

Summarization:

- Discussed about Naïve Bayes algorithm.
- Discussed how to apply Naïve Bayes algorithm to solve the task of Sentiment Analysis in NLP.
- Discussed on an improved baseline to outperform Naïve Bayes.

Question



Thank you!