

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

Text Classification Project



AI VIET NAM
[@aivietnam.edu.vn](https://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, 6 days, 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.

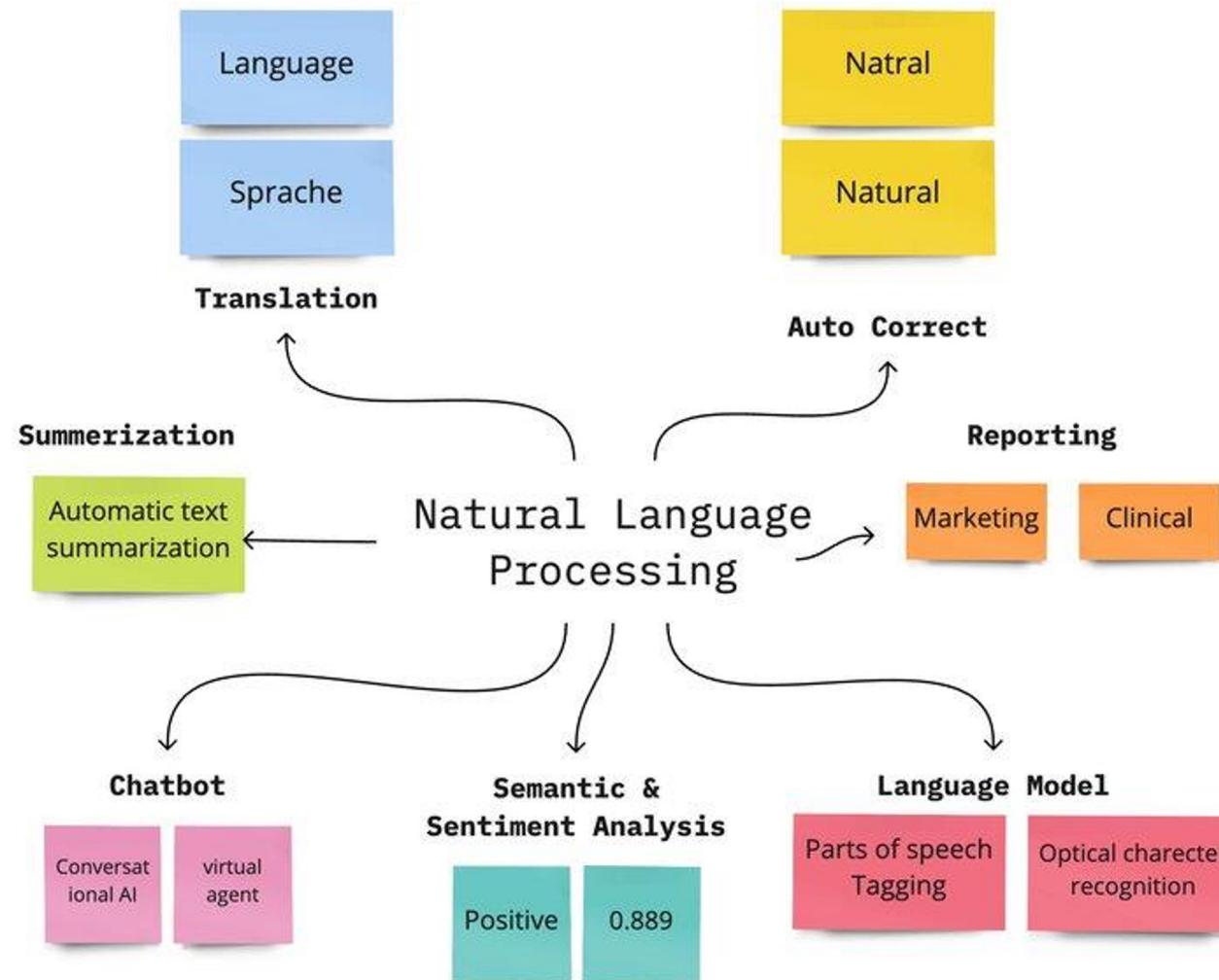
Outline

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

Introduction

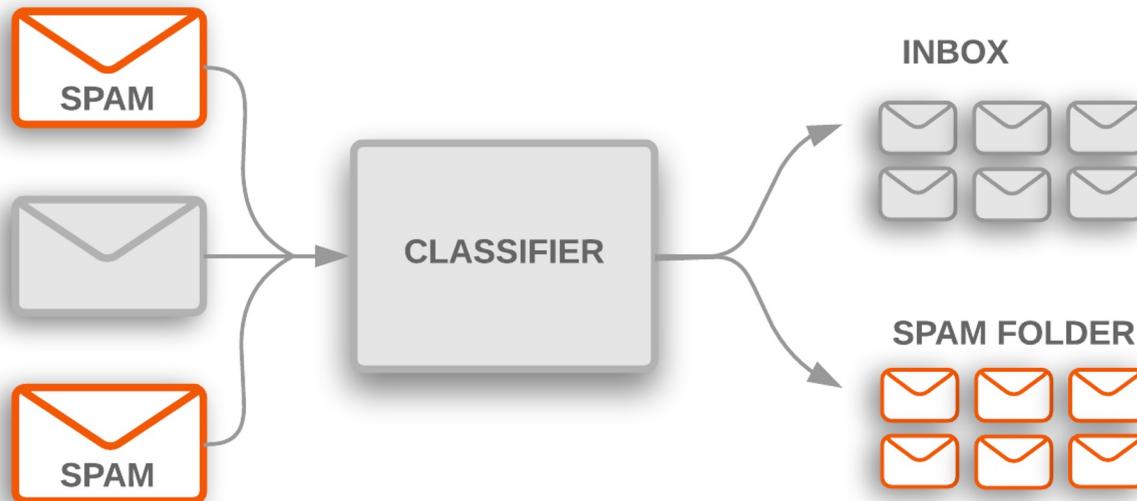
Introduction

❖ Getting Started



Introduction

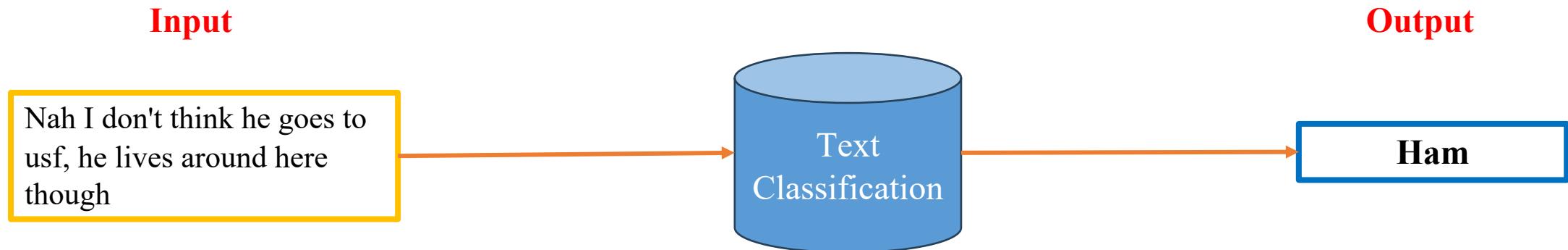
❖ Text Classification



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



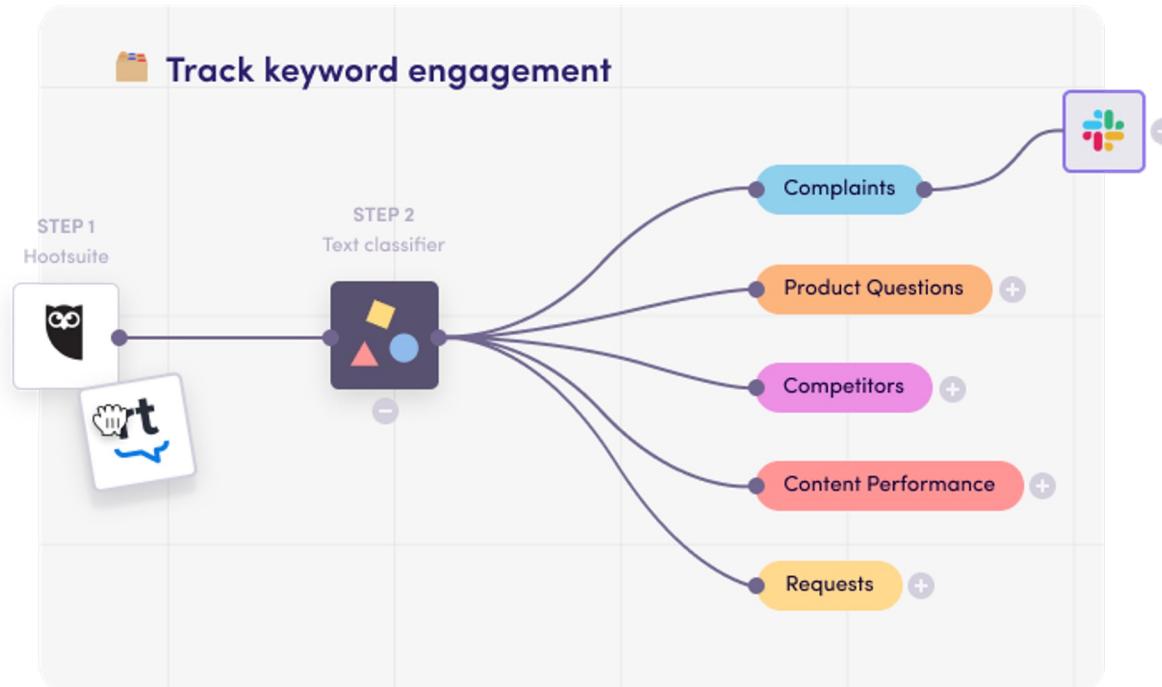
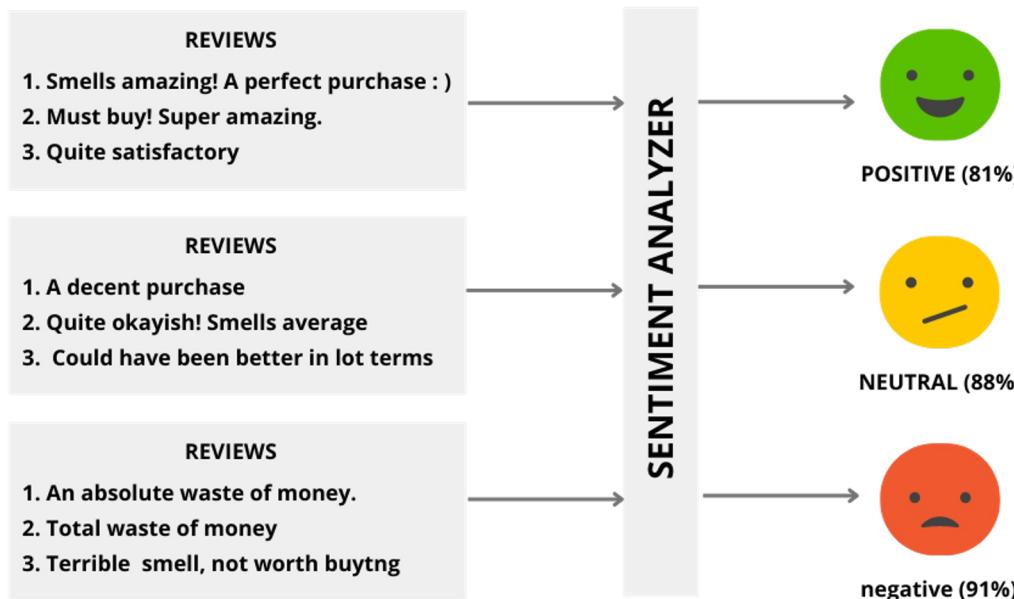
Fragrance-1
(Lavender)



Fragrance-1
(Rose)



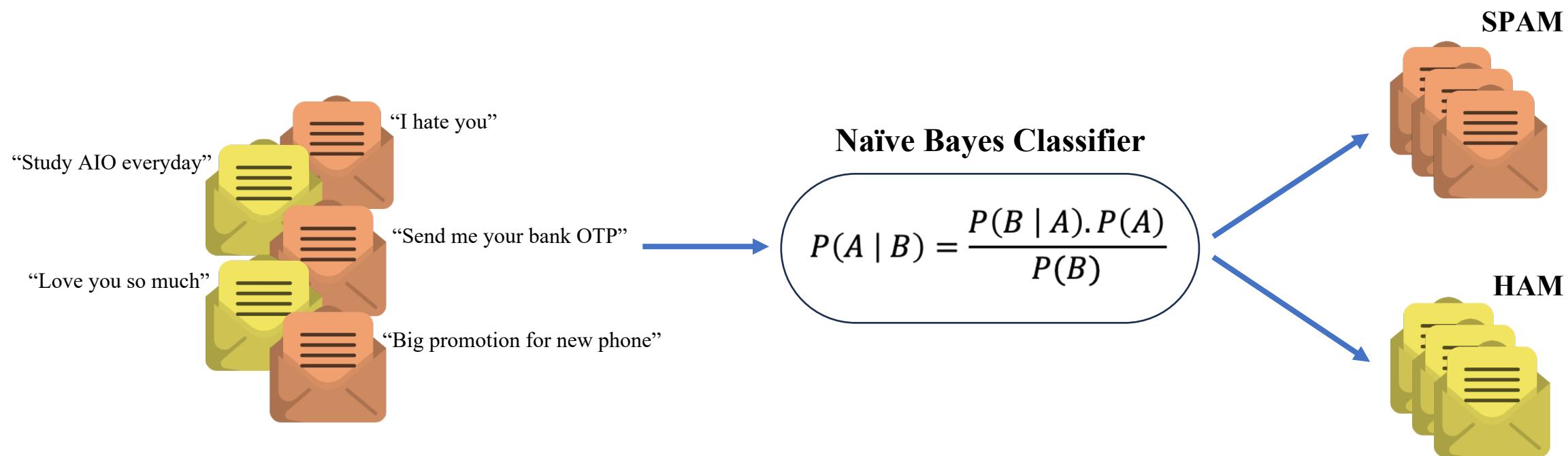
Fragrance-1
(Lemon)



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).



Introduction

❖ Message Classification I/O

Spam	Not spam
 “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”	 “Nah I don't think he goes to usf, he lives around here though”

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

Message Classification

❖ Probability

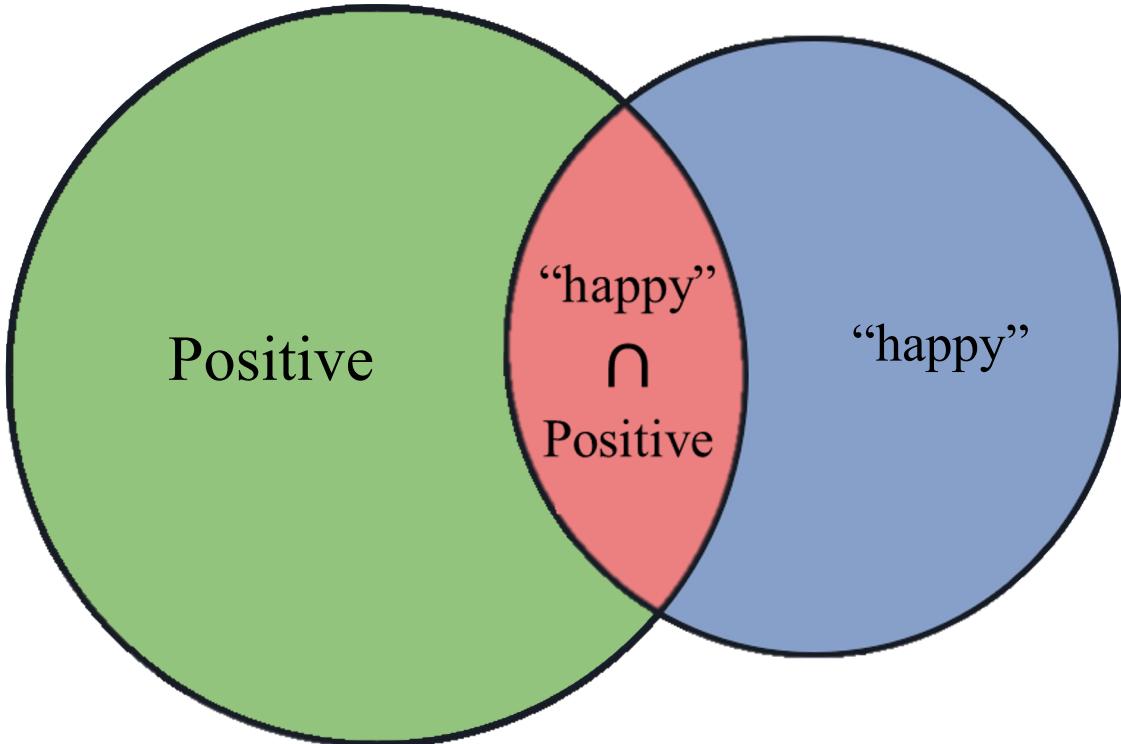
Positive	Positive	Positive	Positive	Positive
Positive	Positive	Positive	<i>Negative</i>	<i>Negative</i>
<i>Negative</i>	<i>Negative</i>	<i>Negative</i>	<i>Negative</i>	<i>Negative</i>
<i>Negative</i>	<i>Negative</i>	<i>Negative</i>	<i>Negative</i>	<i>Negative</i>

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

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

Message Classification

❖ Conditional Probability



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

Message Classification

❖ Bayes' Rule

The content of a message 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



? How do you know which group a word belongs to?



SPAM

"hate", "promotion", "sale", ...



HAM

"happy", "love", "deadline", ...

Message Classification

❖ Bayes' Rule

How do you know which group a word belongs to?



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

$$P(Class \mid Words) = \frac{P(Words \mid Class) \cdot P(Class)}{P(Words)}$$

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

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

!

Message Classification

❖ Word Dependence

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

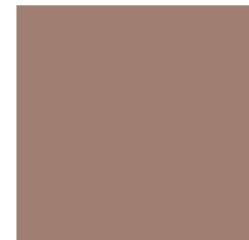
Human:



“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).

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

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

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 Message Classification :

$$P(y | \text{words}) \propto P(y) \cdot \prod_{i=1}^n P(\text{word}_i | y)$$

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	HAM (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 Message 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

❖ Dataset

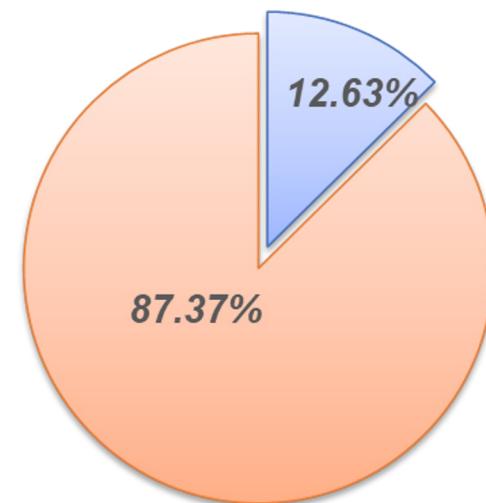
The data set includes 2 columns:

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

▲ Category	▲ Message
ham spam	87% 13%
ham	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

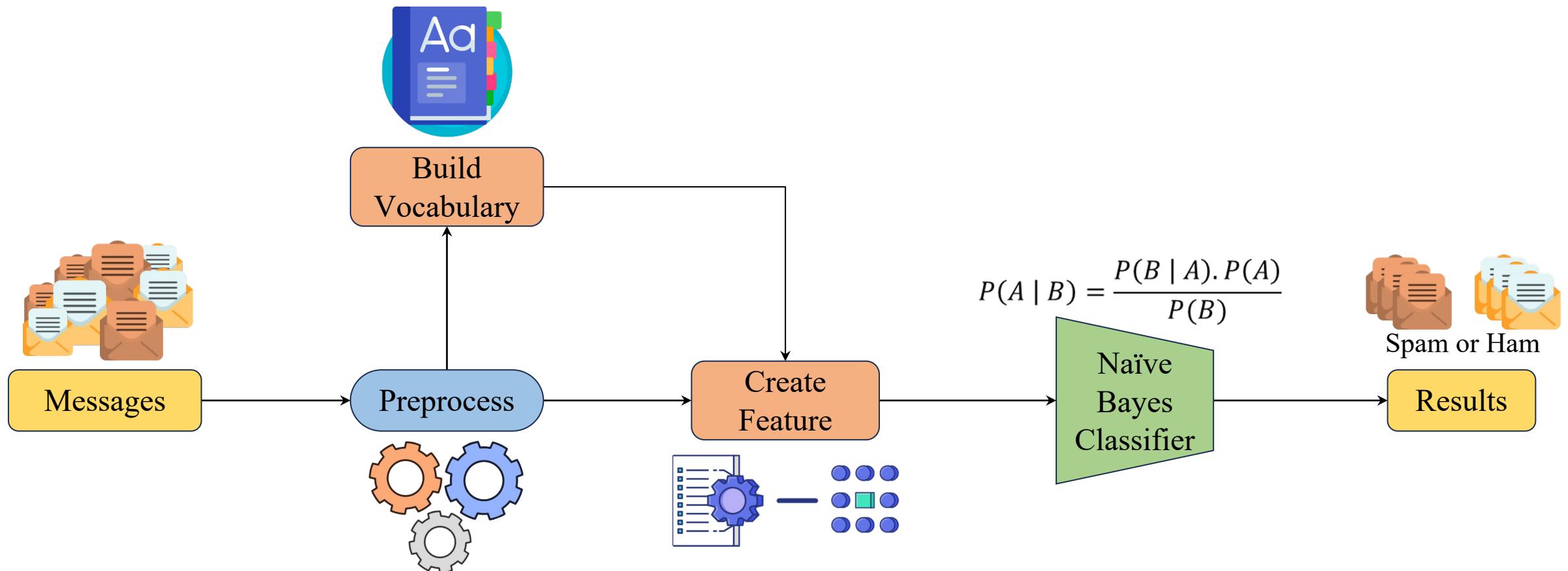
Email Classification

■ Spam ■ Ham



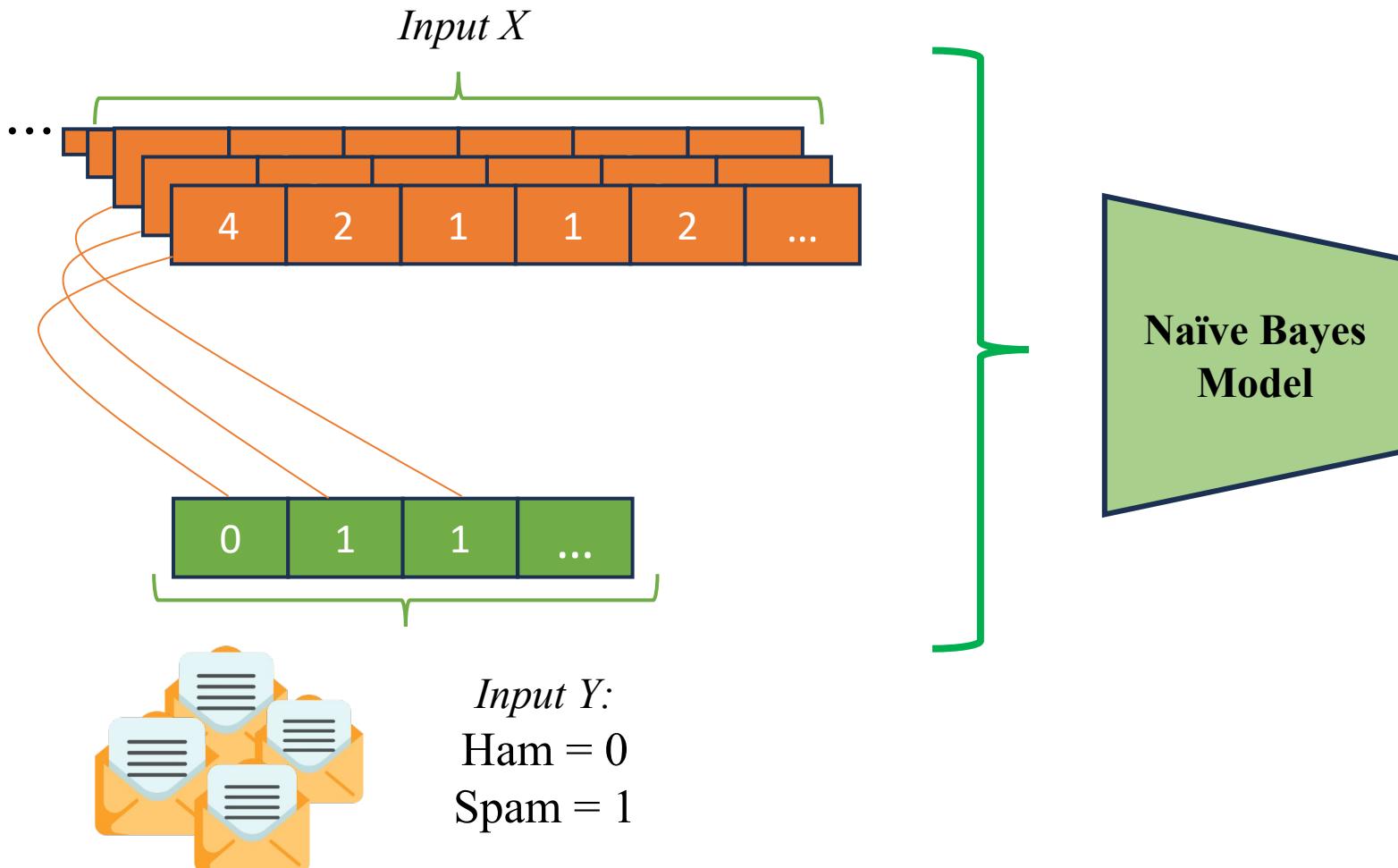
Message Classification

❖ Project Pipeline



Message Classification

❖ Input Naïve Bayes Model

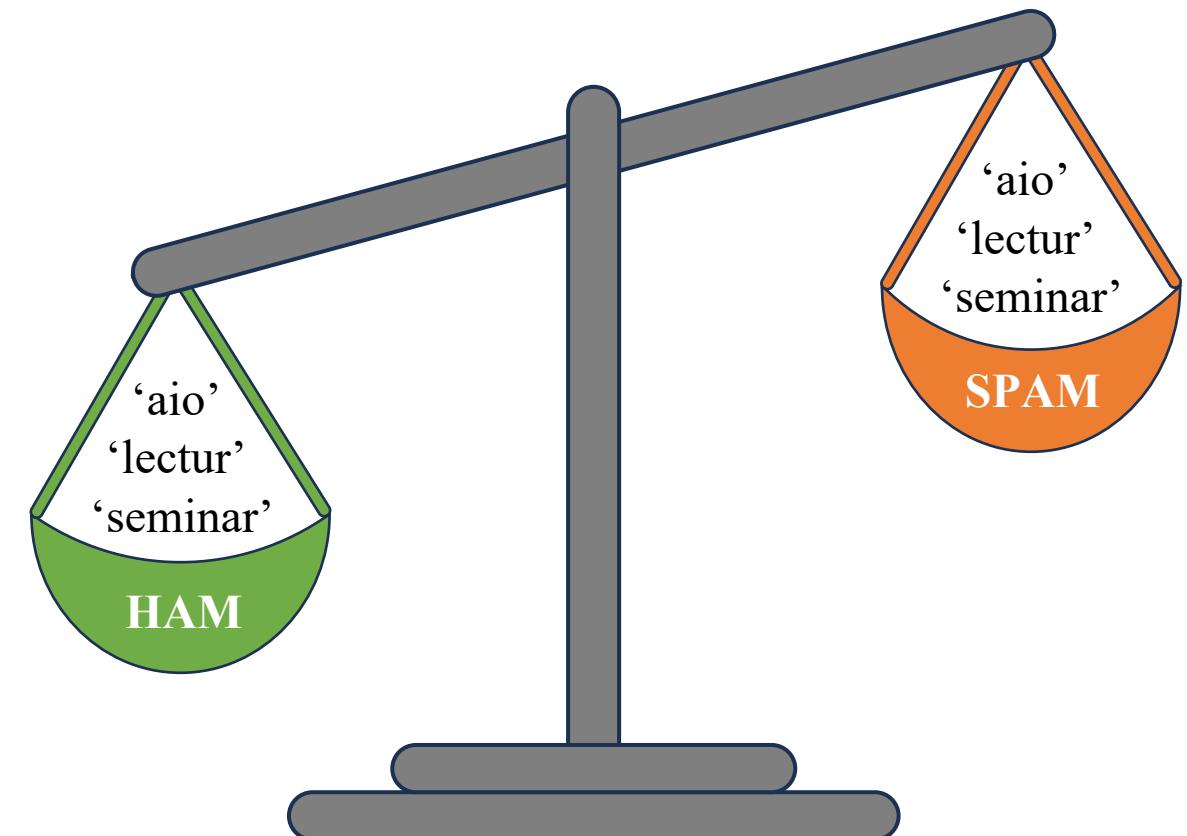


Message Classification

❖ Applying Naïve Bayes with Maximum A Posteriori Estimation

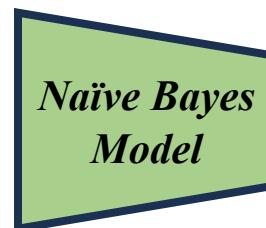
To classification, using MAP to estimation $P(y)$ and $P(word_i | y)$, where y is a label in training set.
Which y has the highest probability is the model result.

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



Message Classification

❖ $P(x_i|y)$



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

aio i studi we love

4 2 1 1 2

0 1 1 0 1

0 3 0 2 0

0 0 1

Suppose Input X:

Suppose Input Y:

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
V	13	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

Message Classification

❖ $P(y)$

*Naïve Bayes
Model*

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

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
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

$$P(Ham) = 0.667$$

$$P(Spam) = 0.334$$

Predict

New message: “I study AIO, I love it”

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

- $\hat{y}_{ham} = (0.667) * (0.231 * 0.154 * 0.308 * 0.231 * 0.231) \approx 0.00038997$
- $\hat{y}_{spam} = (0.334) * (0.6 * 0.1 * 0.1 * 0.6 * 0.1) \approx 0.00012$

$$\hat{y}_{ham} > \hat{y}_{spam} \rightarrow \text{Predict} = \text{Ham}$$

Message Classification

❖ Problem 1 – Zero Probability

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
V	13	5

Diagram illustrating the Naïve Bayes Model calculation:

A blue arrow points from the first row of the left table to the first row of the right table.

Yellow arrows point from the formula $\prod_{i=1}^n P(\text{word}_i | y)$ to the corresponding probability values in the right table.

The right table shows the calculated probabilities for each word category:

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

Message Classification

❖ Problem 1 - Laplacian Smoothing

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

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

Word	Ham	Spam
aio	$(4+1)/(13+5)$	$(0+1)/(5+5)$
i	$(3+1)/(13+5)$	$(3+1)/(5+5)$
studi	$(2+1)/(13+5)$	$(0+1)/(5+5)$
we	$(1+1)/(13+5)$	$(2+1)/(5+5)$
love	$(3+1)/(13+5)$	$(0+1)/(5+5)$
V	13	5

Message Classification

❖ Problem 1 - Laplacian Smoothing

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

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



Word	Ham	Spam
aio	0.278	0.1
i	0.222	0.4
studi	0.167	0.1
we	0.111	0.3
love	0.222	0.1
Sum	1	1

Message Classification

❖ Problem 1 - Laplacian Smoothing

Word	Ham	Spam
aio	0.278	0.1
i	0.222	0.4
studi	0.167	0.1
we	0.111	0.3
love	0.222	0.1
Sum	1	1

Predict

New message: “I study AIO, I love it.”

- $\hat{y}_{ham} = (0.667) * (0.222 * 0.167 * 0.278 * 0.222 * 0.222) \approx 0.0003388$
- $\hat{y}_{spam} = (0.334) * (0.4 * 0.1 * 0.1 * 0.4 * 0.1) \approx 0.00005344$

$\hat{y}_{ham} > \hat{y}_{spam} \rightarrow \text{Predict} = \text{Ham}$

Message Classification

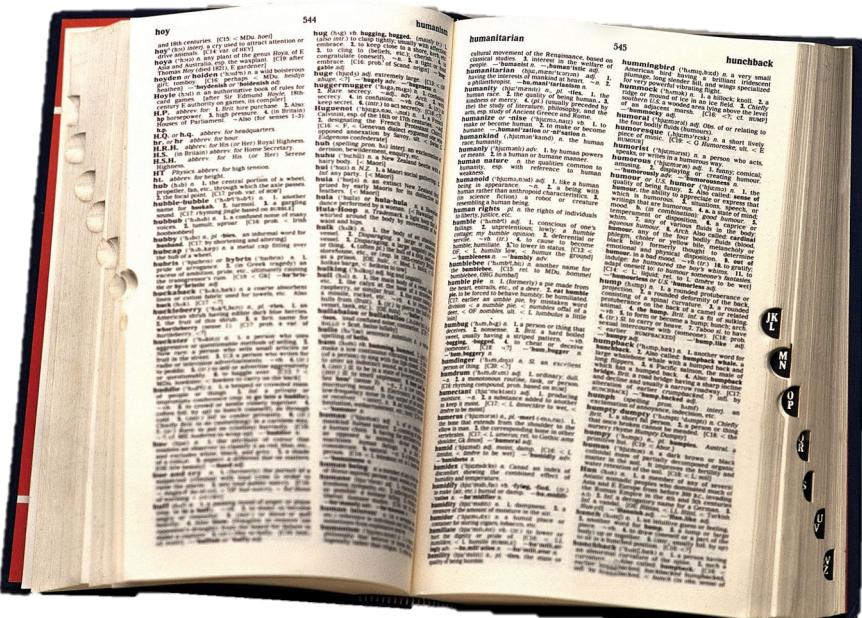
❖ Problem 2 – Log Likelihood

$$\hat{y}_{spam} \approx 0.00005344 \quad \hat{y}_{ham} \approx 0.0003388$$

!

!

In reality, the dictionary is huge



and the message content much longer

Black Friday in July is here! Thư 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 **NBHW**C 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 NBHW (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.

The results will become very very small, gradually approaching 0

Message Classification

❖ Problem 2 – Log Likelihood

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

x	$\text{Log}_{10}(x)$
0.00001	-5
0.001	-3
0	Error
100	3
10000	5

$$\begin{aligned}\hat{y} &= \arg \max_y P(y) \prod_{i=1}^n P(w_i|y) \\ &= \arg \max_y \log P(y) \prod_{i=1}^n P(w_i|y) \quad \text{since } \log(ab) = \log(a) + \log(b) \\ &= \arg \max_y \log P(y) + \log \prod_{i=1}^n P(w_i|y) \\ &= \arg \max_y \log P(y) + \log \sum_{i=1}^n P(w_i|y)\end{aligned}$$

If $a < b$, $\log(a) < \log(b)$
⇒ change from prob to log-prob.
⇒ $(0, 1)$ to $(-\infty, 0)$

Message Classification

❖ Problem 2 – Log Likelihood

Predict

New message: “I study AIO, I love it.”

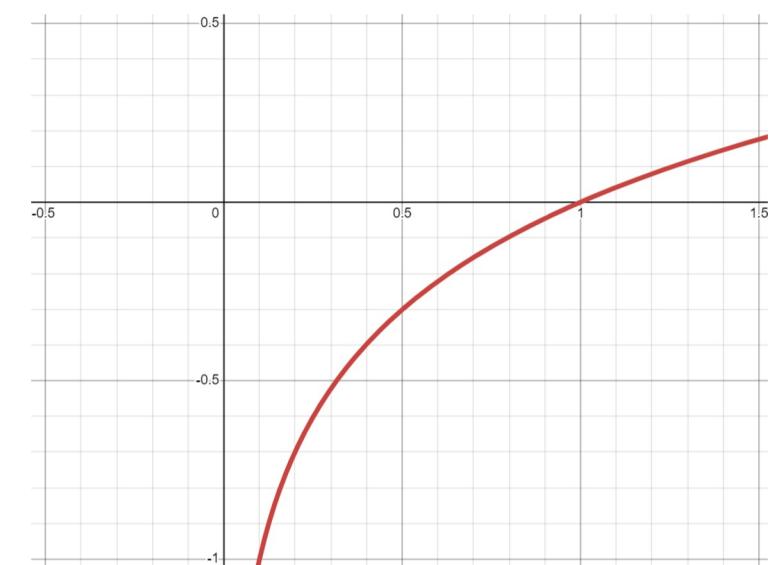
$$\begin{aligned}\hat{y}_{ham} &= [\log(0.222) + \log(0.167) + \log(0.278) + \log(0.222) + \log(0.222)] + \log(0.667) \\ &\approx -3.294 - 0.175 \approx -3.469\end{aligned}$$

$$\begin{aligned}\hat{y}_{spam} &= [\log(0.4) + \log(0.1) + \log(0.1) + \log(0.4) + \log(0.1)] + \log(0.334) \\ &\approx -3.796 - 0.476 \approx -4.272\end{aligned}$$

Since $a < b$, $\log(a) < \log(b) \Rightarrow$ A smaller number (more negative) corresponds to a smaller probability, hence:

$$\hat{y}_{ham} > \hat{y}_{spam} \rightarrow \text{Predict} = \text{Ham}$$

Word	Ham	Spam
aio	0.278	0.1
i	0.222	0.4
studi	0.167	0.1
we	0.111	0.3
love	0.222	0.1



Break - Quiz

Break - Quiz

Câu 1: Mục tiêu chính của bài toán Text Classification là gì?

- A. Tóm tắt nội dung văn bản.
- B. Dịch văn bản sang ngôn ngữ khác.
- C. Phân loại văn bản vào các nhóm định trước.
- D. Tạo ra văn bản mới dựa trên văn bản đã cho.

Câu 2: Trong bài toán phân loại email, hai lớp chính thường được sử dụng là gì?

- A. Quan trọng và Không quan trọng.
- B. Spam và Ham (không phải spam).
- C. Cá nhân và Công việc.
- D. Đã đọc và Chưa đọc.

Câu 3: Giả định "Naive" trong thuật toán Naive Bayes là gì?

- A. Tất cả email đều có xác suất như nhau để là Spam.
- B. Độ dài của email quyết định nó là Spam hay Ham.
- C. Sự xuất hiện của một từ không ảnh hưởng đến sự xuất hiện của các từ khác.
- D. Các email Spam luôn chứa các từ khóa cụ thể.

Câu 4: Khi áp dụng Naive Bayes cho phân loại email, tại sao có thể bỏ qua mẫu số $P(w)$ trong công thức Bayes?

- A. Vì $P(w)$ luôn bằng 1.
- B. Vì $P(w)$ không ảnh hưởng đến việc so sánh xác suất giữa các lớp.
- C. Vì $P(w)$ chỉ quan trọng trong trường hợp có nhiều hơn 2 lớp.
- D. Vì $P(w)$ luôn nhỏ hơn $P(c|w)$.

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).

