



Naive Bayes clearly explained!

Text classification using naïve bayes

Conditional Probability

1. In probability theory, conditional probability is a measure of the probability of an event given that another event has already occurred.
2. If the event of interest is A and the event B is assumed to have occurred, "the conditional probability of A given B ", or "the probability of A under the condition B ", is usually written as $P(A|B)$, or sometimes $P_B(A)$.

Example

Chances of cough

The probability that any given person has a cough on any given day maybe only 5%. But if we know or assume that the person has a cold, then they are much more likely to be coughing. The conditional probability of coughing given that person have a cold might be a much higher 75%.

Marbles in a bag

2 blue and 3 red marbles are in a bag.

What are the chances of getting a blue marble?

???

Marbles in a bag

2 blue and 3 red marbles are in a bag.

What are the chances of getting a blue marble?

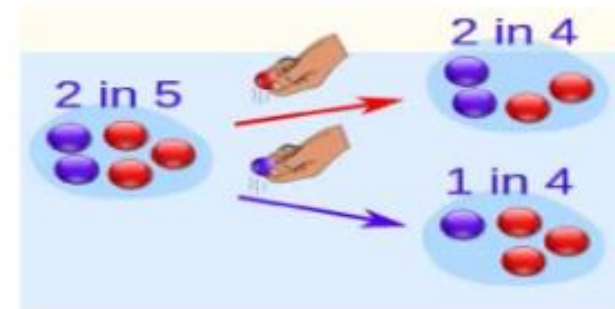
Answer: - The chance is 2 in 5

Marbles in a bag

But after taking one out of these chances,
situation may change!

So the next time:

1. if we got a red marble before, then the chance of a blue marble next is 2 in 4
2. if we got a blue marble before, then the chance of a blue marble next is 1 in 4



Probability given pre-condition

Likewise:-

Drawing a second ace from a deck given we got the first ace

Finding the probability of having a disease given you were tested positive

Finding the probability of liking Harry Potter given we know the person likes fiction.

Bayes Theorem

1. In probability theory and statistics, Bayes' theorem (alternatively Bayes' law or Bayes' rule) describes the probability of an event, based on prior knowledge of conditions that might be related to the event.
2. For example, if cancer is related to age, then, using Bayes' theorem, a person's age can be used to more accurately to assess the probability that they have cancer, compared to the assessment of the probability of cancer made without knowledge of the person's age.

Bayes Theorem

Likelihood of the Evidence
given that the Hypothesis is
True

Prior Probability of the
Hypothesis

$$P(H \setminus E) = \frac{P(E \setminus H) * P(H)}{P(E)}$$

Prior probability of the
Hypothesis given that the
Evidence is True

Prior probability that
the evidence is True

Bayes Theorem

The Formula for Bayes' theorem

$$P(H | E) = \frac{P(E | H) * P(H)}{P(E)}$$

where

1. $P(H)$ is the probability of hypothesis H being true. This is known as the prior probability.
2. $P(E)$ is the probability of the evidence (regardless of the hypothesis).
3. $P(E|H)$ is the probability of the evidence given that hypothesis is true.
4. $P(H|E)$ is the probability of the hypothesis given that the evidence is there.

Naïve Bayes

- ❖ 3 types
- ❖ Gaussian Naïve Bayes
- ❖ Binomial Naïve Bayes
- ❖ Multinomial Naïve Bayes

Email Classification



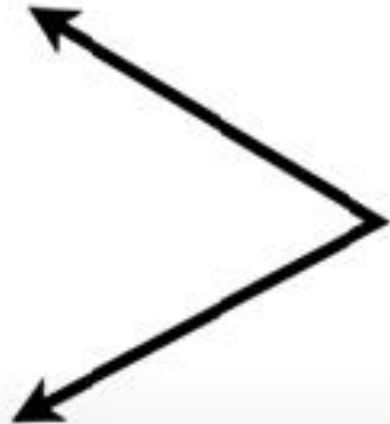
Imagine we received
normal messages from
friends and family...

Email Classification



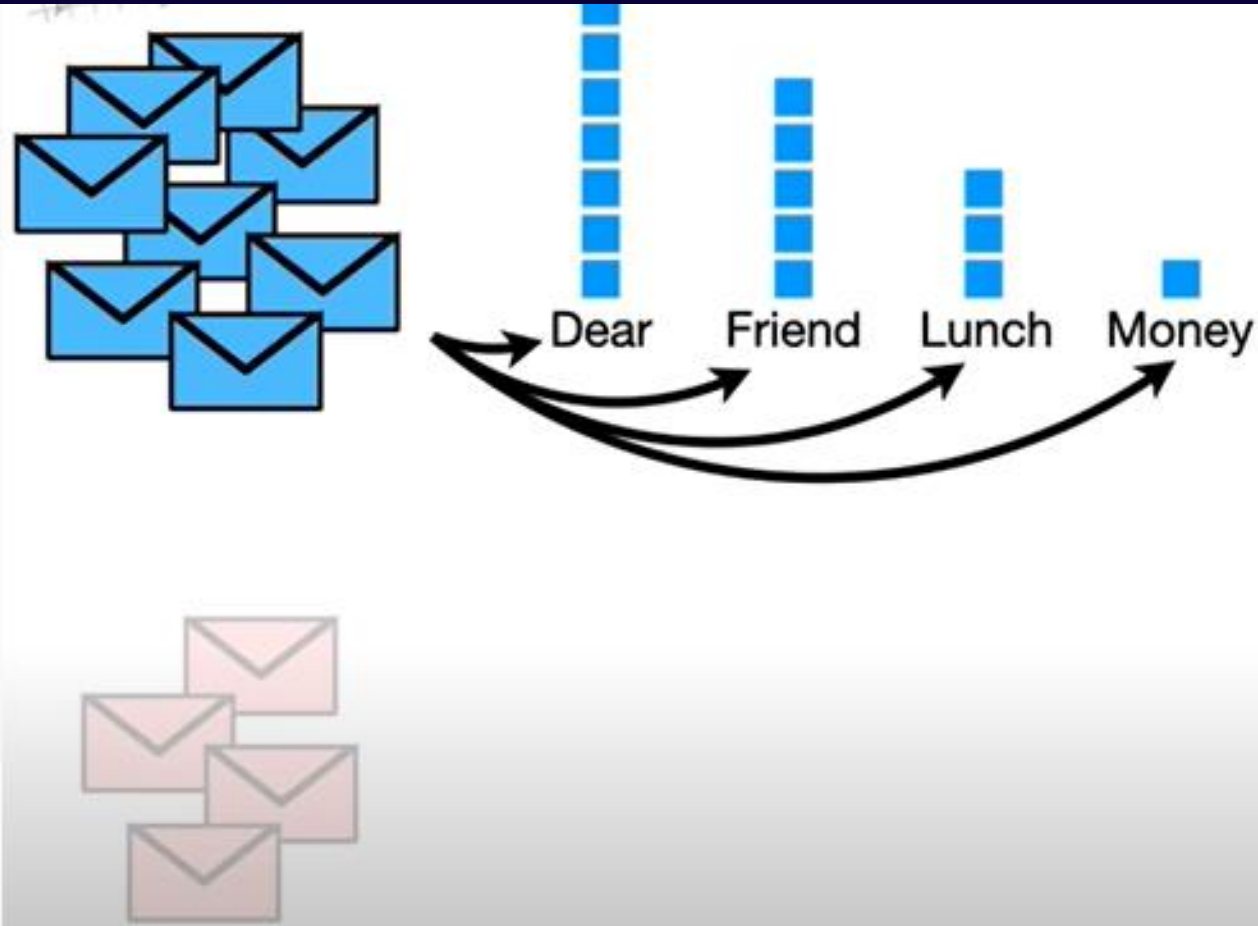
...and we also received
spam (unwanted
messages that are usually
scams or unsolicited
advertisements)...

Email Classification



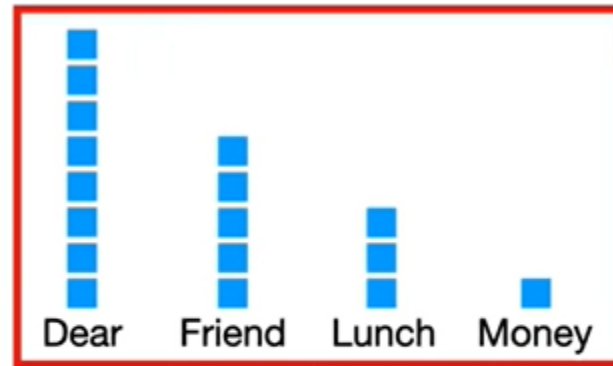
...and we wanted to filter out the **spam** messages.

Email Classification



So, the first thing we do is make a **histogram** of all the words that occur in the **normal messages** from friends and family.

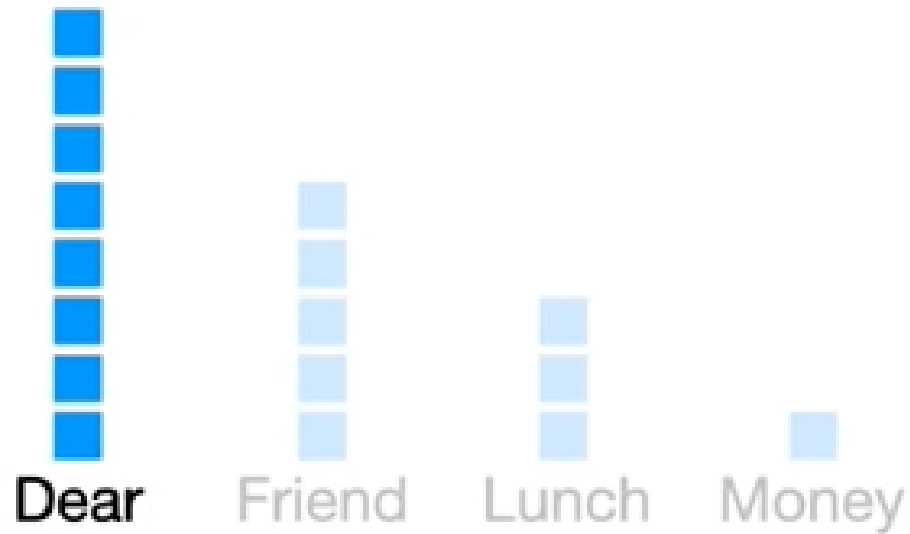
Email Classification



...divided by **17**, the total number of words in all of the **normal messages**.

$$p(\text{Dear} \mid \text{Normal}) = \frac{8}{17}$$

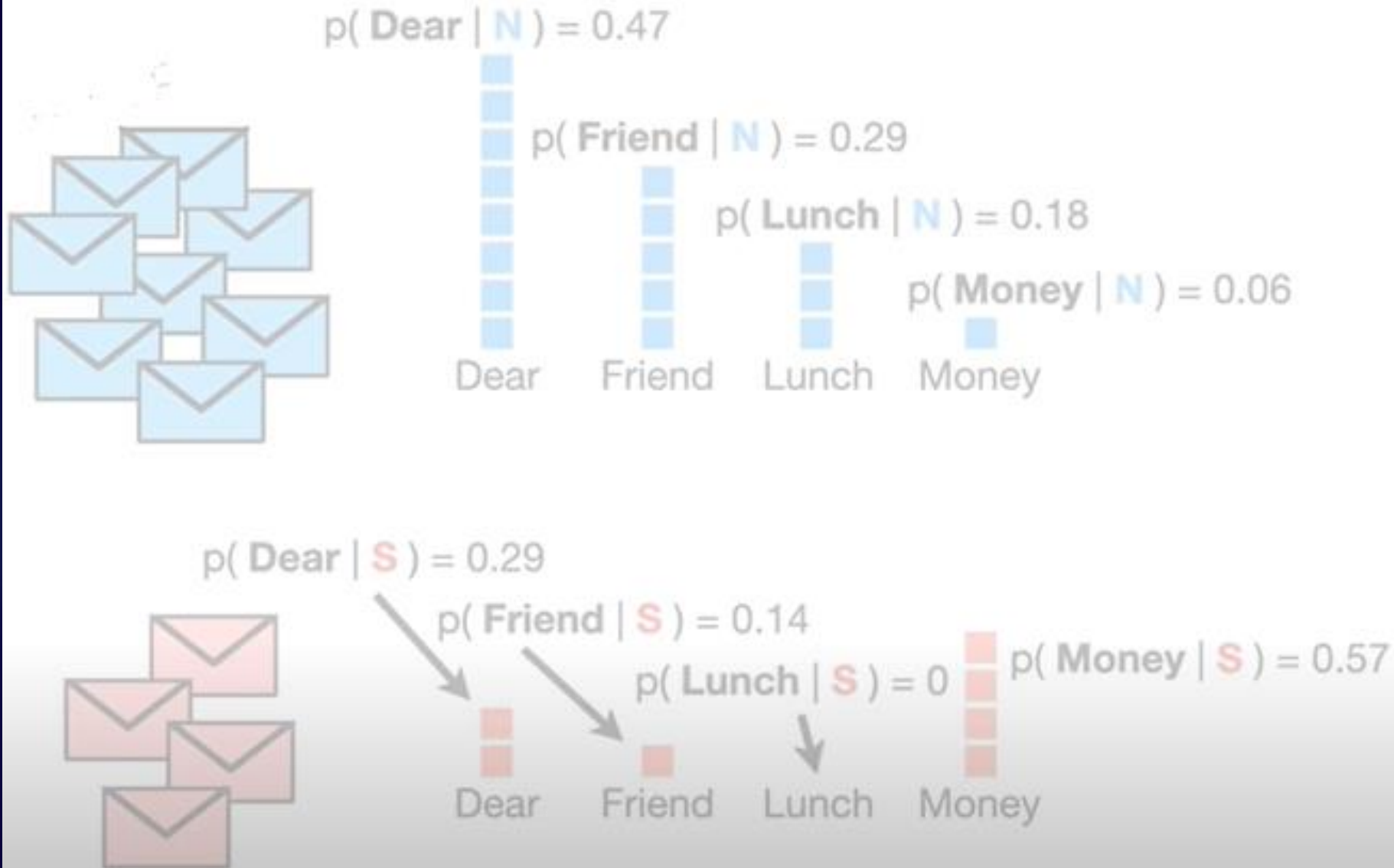

Email Classification



And that gives us **0.47**.

$$p(\text{Dear} \mid \text{Normal}) = \frac{8}{17} = 0.47$$

Email Classification



Bam!

Email Classification



$$\begin{aligned}p(\text{Dear} \mid \mathbf{N}) &= 0.47 \\p(\text{Friend} \mid \mathbf{N}) &= 0.29 \\p(\text{Lunch} \mid \mathbf{N}) &= 0.18 \\p(\text{Money} \mid \mathbf{N}) &= 0.06\end{aligned}$$



$$\begin{aligned}p(\text{Dear} \mid \mathbf{S}) &= 0.29 \\p(\text{Friend} \mid \mathbf{S}) &= 0.14 \\p(\text{Lunch} \mid \mathbf{S}) &= 0.00 \\p(\text{Money} \mid \mathbf{S}) &= 0.57\end{aligned}$$

Terminology Alert!!!

Because we have calculated the probabilities of discrete, individual words, and not the probability of something continuous, like weight or height, these **Probabilities** are also called **Likelihoods**.



Email Classification



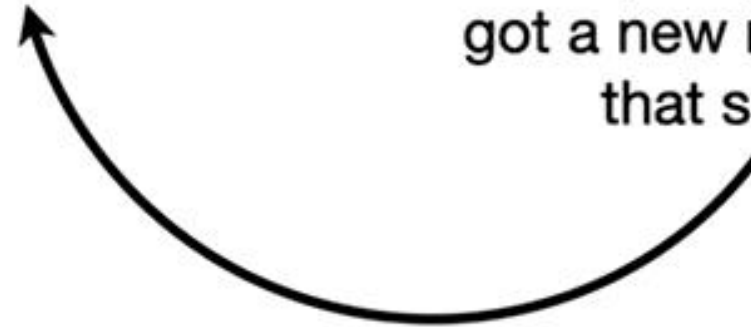
$$\begin{aligned}p(\text{Dear} \mid \mathbf{N}) &= 0.47 \\p(\text{Friend} \mid \mathbf{N}) &= 0.29 \\p(\text{Lunch} \mid \mathbf{N}) &= 0.18 \\p(\text{Money} \mid \mathbf{N}) &= 0.06\end{aligned}$$



$$\begin{aligned}p(\text{Dear} \mid \mathbf{S}) &= 0.29 \\p(\text{Friend} \mid \mathbf{S}) &= 0.14 \\p(\text{Lunch} \mid \mathbf{S}) &= 0.00 \\p(\text{Money} \mid \mathbf{S}) &= 0.57\end{aligned}$$

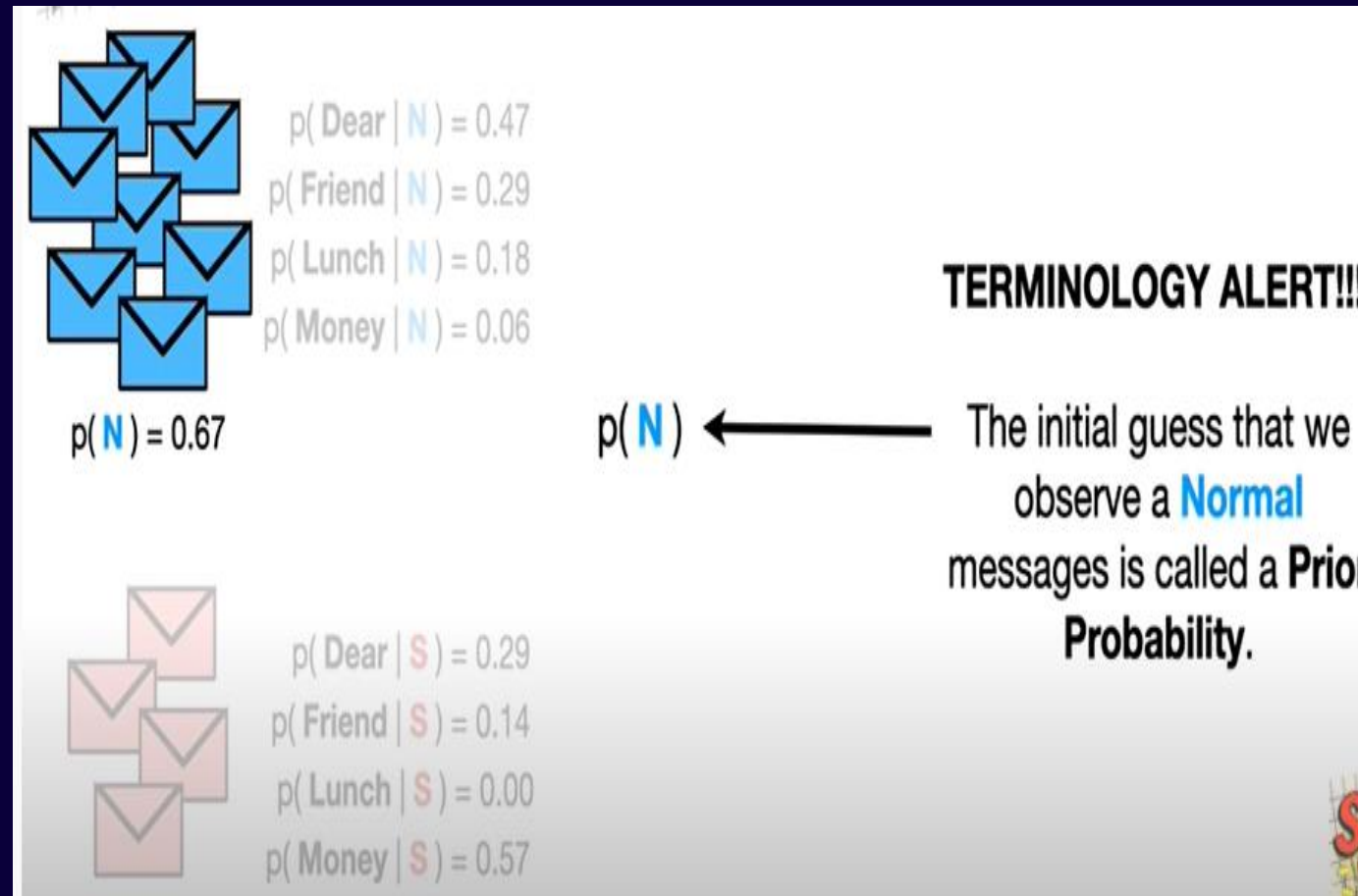
Dear Friend

Now, imagine we
got a new message
that said:

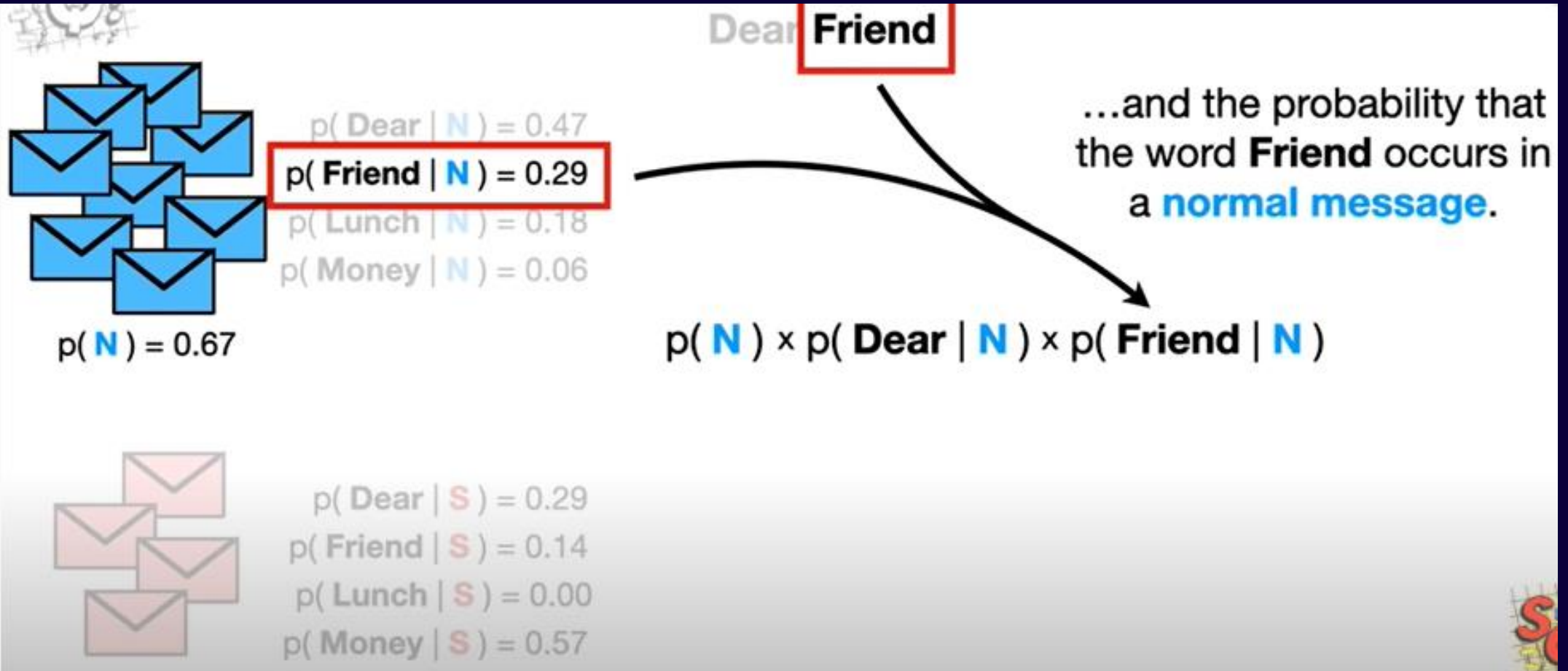


Email Classification

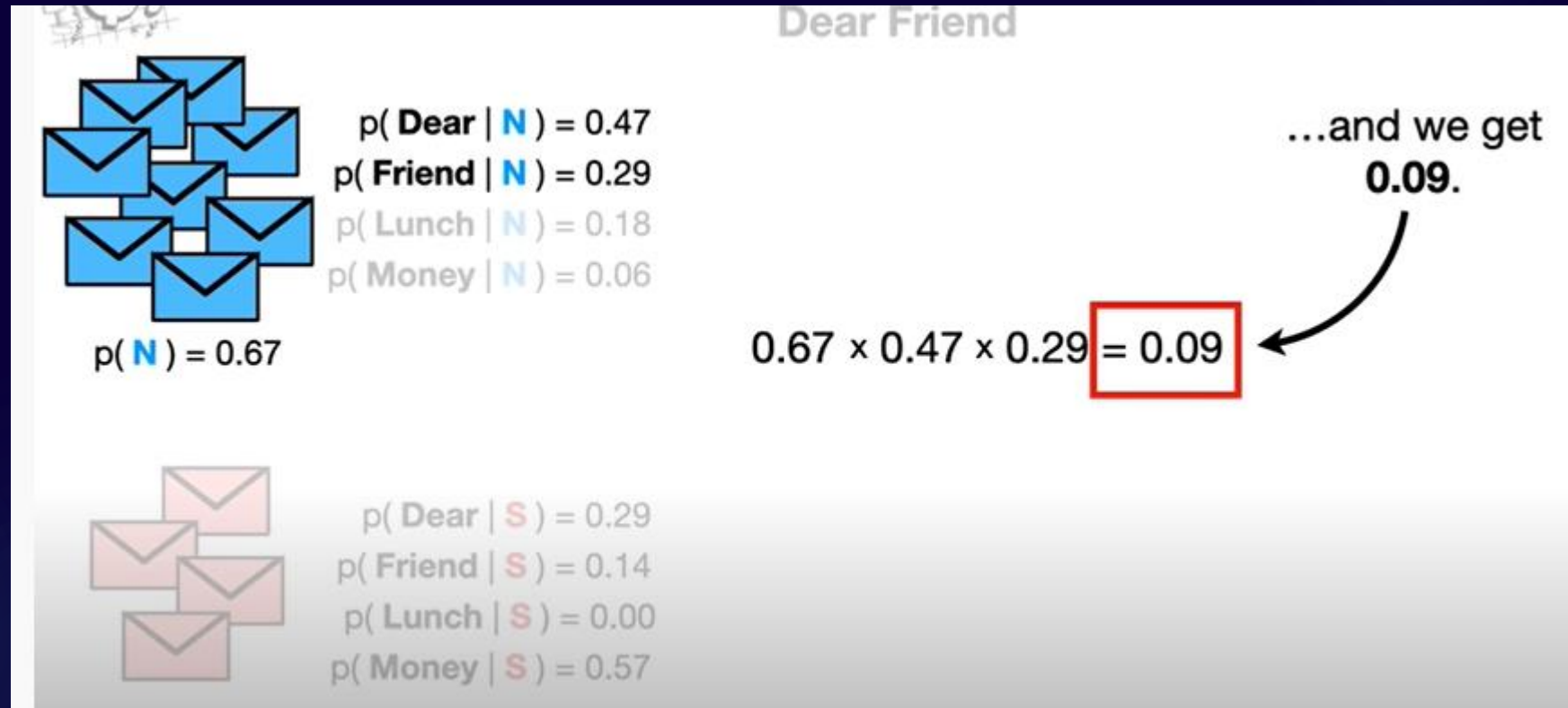
- ❖ Start with an initial guess of the probability it is a normal message
- ❖ Let 8 normal messages
- ❖ 12 total messages
- ❖ This $P(N) = 8/12 = 0.67$



Email Classification



Email Classification



$P(S) = \text{no. of spam messages} / \text{total messages} = 4/12 = 0.33$


Think of .01 as the score that Dear Friend gets if it is **Spam**

$$P(S | \text{Dear Friend}) = P(S) * P(\text{Dear} | S) * P(\text{Friend} | S)$$


$$= 0.33 * 0.29 * 0.14 = 0.01$$

Email Classification

$p(\mathbf{N} \mid \text{Dear Friend}) \propto 0.09$


 $p(\mathbf{N}) = 0.67$

$p(\text{Dear} \mid \mathbf{N}) = 0.47$
 $p(\text{Friend} \mid \mathbf{N}) = 0.29$
 $p(\text{Lunch} \mid \mathbf{N}) = 0.18$
 $p(\text{Money} \mid \mathbf{N}) = 0.06$


Dear Friend 

...we will decide that **Dear Friend** is a **Normal Message**

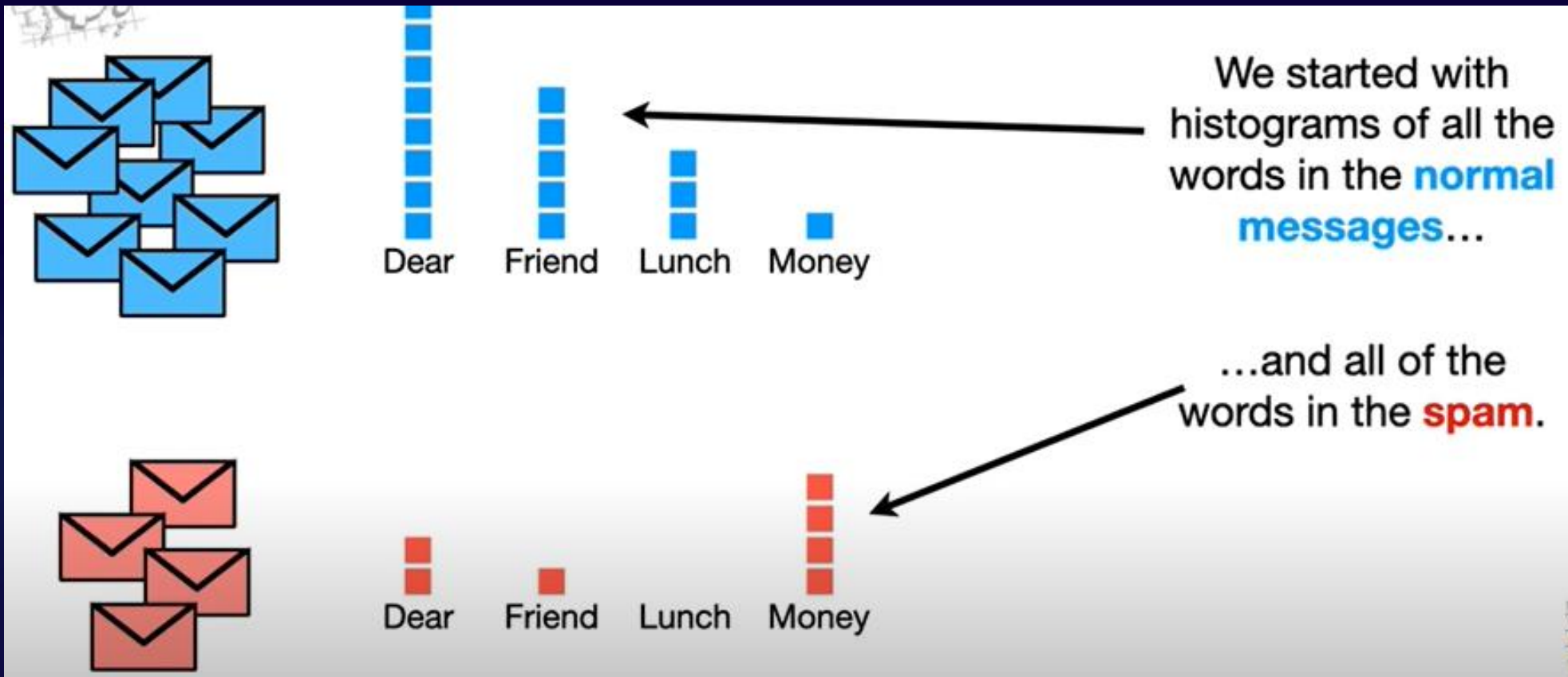
$0.33 \times 0.29 \times 0.14 = 0.01 \propto p(\mathbf{S} \mid \text{Dear Friend})$

 $p(\mathbf{S}) = 0.33$

$p(\text{Dear} \mid \mathbf{S}) = 0.29$
 $p(\text{Friend} \mid \mathbf{S}) = 0.14$
 $p(\text{Lunch} \mid \mathbf{S}) = 0.00$
 $p(\text{Money} \mid \mathbf{S}) = 0.57$



Recap



Recap

- ❖ We started with histogram of words in normal message and spam
- ❖ Calculate probability of each word in normal message or spam
- ❖ Calculate probability of seeing a normal message
- ❖ This is based on training dataset
- ❖ Calculate probability of seeing a spam

Recap

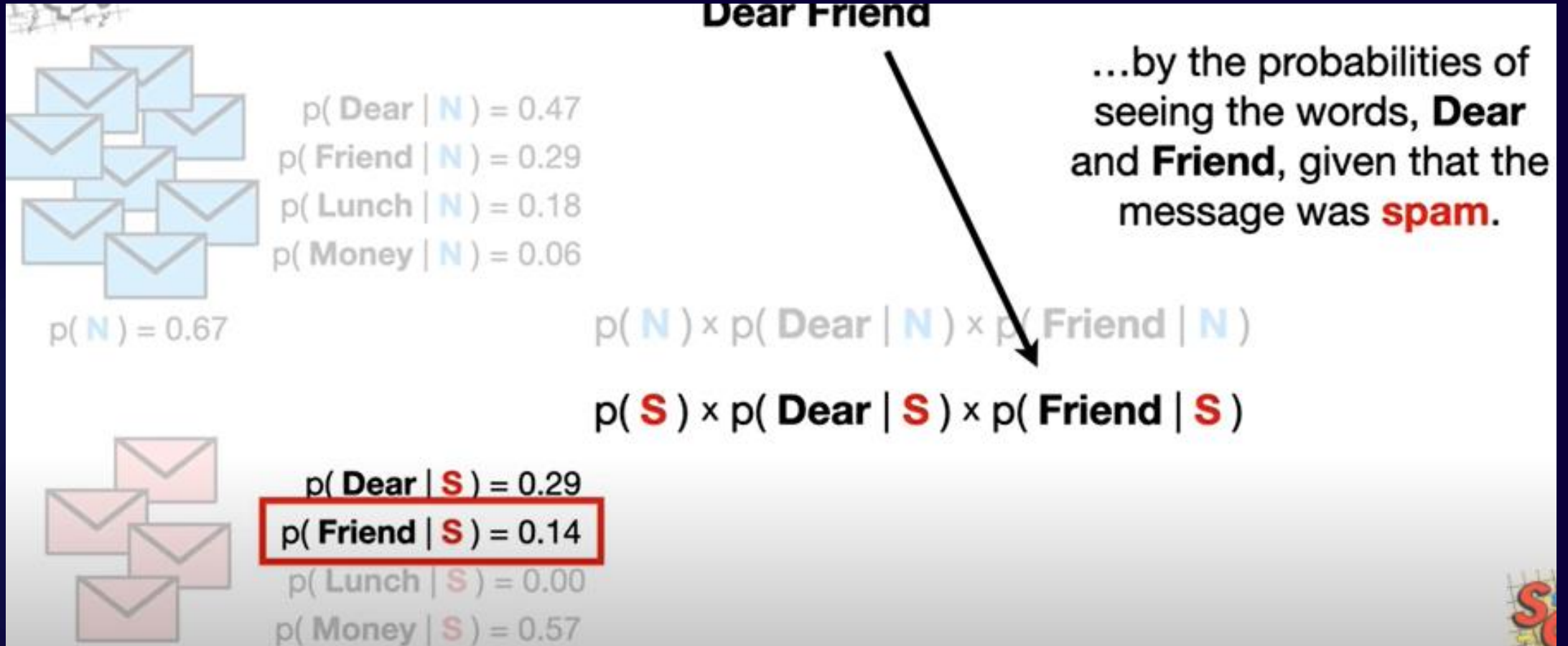
Dear Friend

...by the probabilities of seeing the words, **Dear** and **Friend**, given that the message was **spam**.

$p(\text{Dear} | \text{N}) = 0.47$
 $p(\text{Friend} | \text{N}) = 0.29$
 $p(\text{Lunch} | \text{N}) = 0.18$
 $p(\text{Money} | \text{N}) = 0.06$
 $p(\text{N}) = 0.67$

$p(\text{Dear} | \text{S}) = 0.29$
 $p(\text{Friend} | \text{S}) = 0.14$
 $p(\text{Lunch} | \text{S}) = 0.00$
 $p(\text{Money} | \text{S}) = 0.57$

$p(\text{N}) \times p(\text{Dear} | \text{N}) \times p(\text{Friend} | \text{N})$
 $p(\text{S}) \times p(\text{Dear} | \text{S}) \times p(\text{Friend} | \text{S})$



Recap



$$p(\mathbf{N}) = 0.67$$

$$\begin{aligned}p(\text{Dear} \mid \mathbf{N}) &= 0.47 \\p(\text{Friend} \mid \mathbf{N}) &= 0.29 \\p(\text{Lunch} \mid \mathbf{N}) &= 0.18 \\p(\text{Money} \mid \mathbf{N}) &= 0.06\end{aligned}$$

Dear Friend



Now that we understand the basics of how **Naive Bayes Classification** works...

$$p(\mathbf{N}) \times p(\text{Dear} \mid \mathbf{N}) \times p(\text{Friend} \mid \mathbf{N}) = 0.09$$

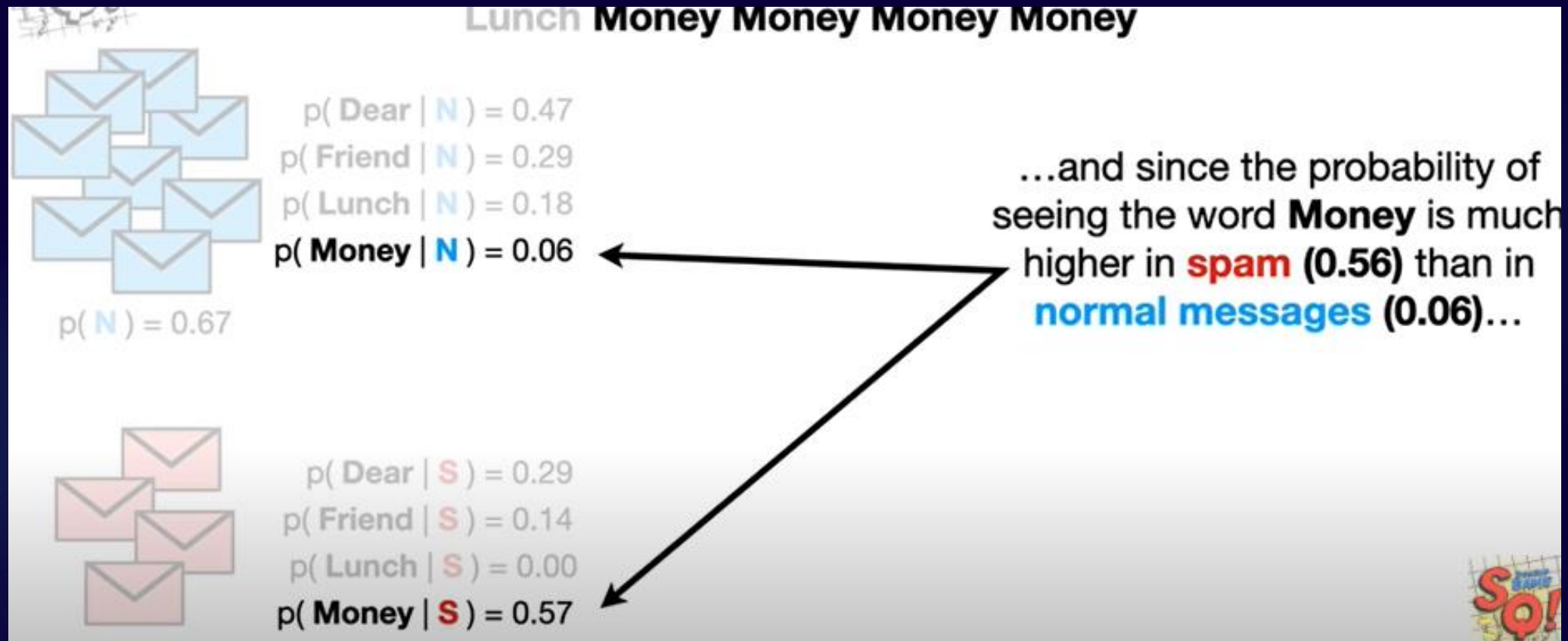
$$p(\mathbf{S}) \times p(\text{Dear} \mid \mathbf{S}) \times p(\text{Friend} \mid \mathbf{S}) = 0.01$$



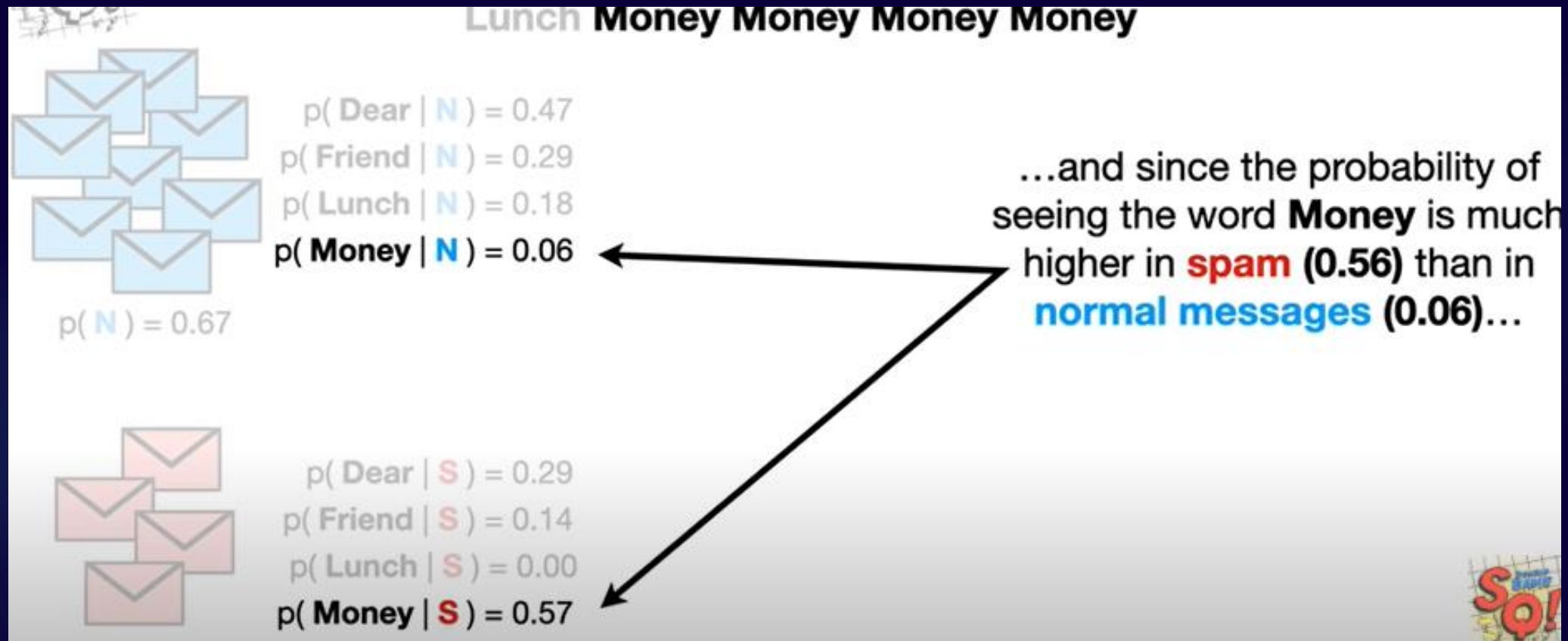
$$\begin{aligned}p(\text{Dear} \mid \mathbf{S}) &= 0.29 \\p(\text{Friend} \mid \mathbf{S}) &= 0.14 \\p(\text{Lunch} \mid \mathbf{S}) &= 0.00 \\p(\text{Money} \mid \mathbf{S}) &= 0.57\end{aligned}$$



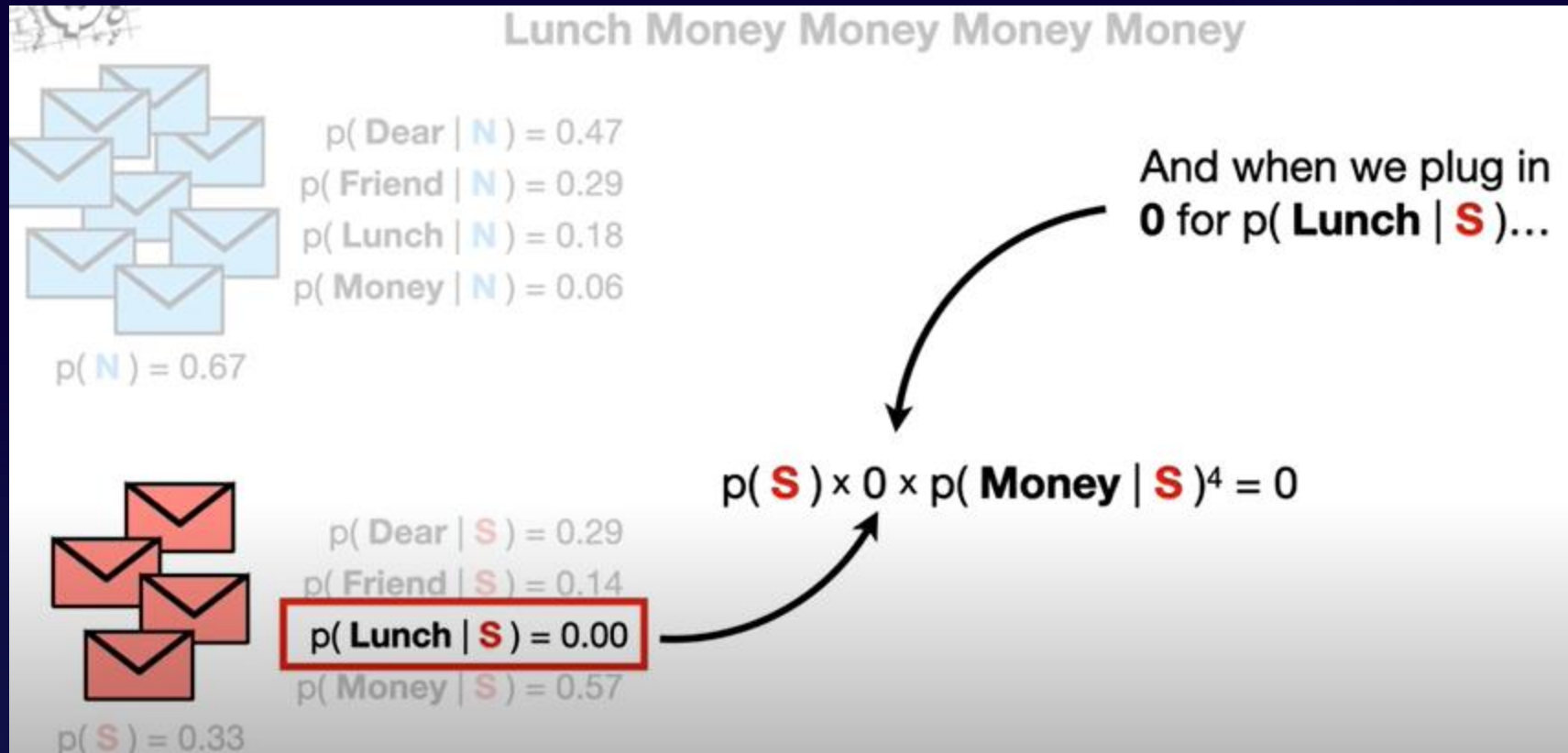
Question



Question

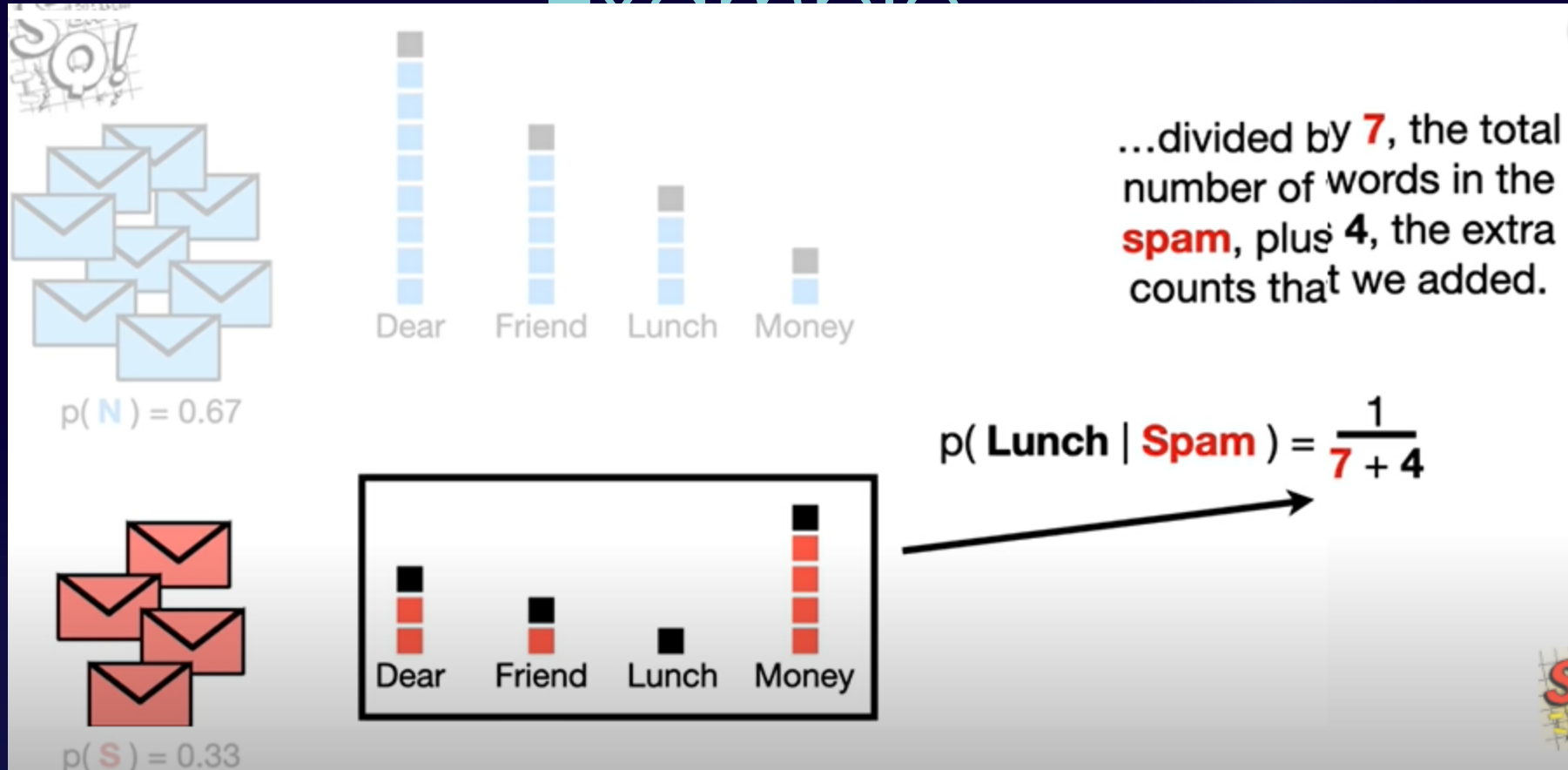


Example



$P(S) = 0$ here since we have not yet seen the word lunch yet in spam

Example



We add 1/any number to all words
Now $P(\text{Lunch} | \text{spam})$ becomes non zero

Example


Lunch Money Money Money Money

...but now when we calculate the value for **spam**, we get a value > 0

↓

$p(\mathbf{N}) \times p(\mathbf{Lunch} | \mathbf{N}) \times p(\mathbf{Money} | \mathbf{N})^4 = 0.00001$

$p(\mathbf{S}) \times p(\mathbf{Lunch} | \mathbf{S}) \times p(\mathbf{Money} | \mathbf{S})^4 = 0.00122$



Normal (N) Probabilities:

- $p(\mathbf{N}) = 0.67$
- $p(\mathbf{Dear} | \mathbf{N}) = 0.43$
- $p(\mathbf{Friend} | \mathbf{N}) = 0.29$
- $p(\mathbf{Lunch} | \mathbf{N}) = 0.19$
- $p(\mathbf{Money} | \mathbf{N}) = 0.10$

Spam (S) Probabilities:

- $p(\mathbf{S}) = 0.33$
- $p(\mathbf{Dear} | \mathbf{S}) = 0.27$
- $p(\mathbf{Friend} | \mathbf{S}) = 0.18$
- $p(\mathbf{Lunch} | \mathbf{S}) = 0.09$
- $p(\mathbf{Money} | \mathbf{S}) = 0.45$


Now we calculate the probability of message being spam/normal
We get $P(S) > P(N)$
So we say that message Lunch Money.... Is **SPAM**

SPAM!!!



Why Naïve Bayes is naive

Why Naïve Bayes is Naive



$p(\mathbf{N}) = 0.67$

$p(\mathbf{Dear} | \mathbf{N}) = 0.43$
 $p(\mathbf{Friend} | \mathbf{N}) = 0.29$
 $p(\mathbf{Lunch} | \mathbf{N}) = 0.19$
 $p(\mathbf{Money} | \mathbf{N}) = 0.10$

Score for **Dear Friend** =
 $p(\mathbf{N}) \times p(\mathbf{Dear} | \mathbf{N}) \times p(\mathbf{Friend} | \mathbf{N}) = 0.08$

Score for **Friend Dear** =
 $p(\mathbf{N}) \times p(\mathbf{Friend} | \mathbf{N}) \times p(\mathbf{Dear} | \mathbf{N}) = 0.08$

In other words, regardless of how the words are ordered, we get **0.08**.

- ❖ Treats all word orders the same
- ❖ Regardless of order of word dear friend/ friend dear
- ❖ Probability is 0.08

Why Naïve Bayes is Naive

- ❖ Treating all word orders equal is very different from how we communicate
- ❖ Every language has grammar rules
- ❖ And common phrases but Naïve Bayes ignores all that
- ❖ Naïve Bayes treats language like a bag full of words
- ❖ We can say that Naïve Bayes has a high bias
- ❖ Has low variance

Theory Behind Naïve Bayes

Theory Behind Naïve Baiyes

Bayes' Theorem

Bayes' Theorem finds the probability of an event occurring given the probability of another event that has already occurred. Bayes' theorem is stated mathematically as the following equation:

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

Theory Behind Naïve Baiyes

Naive assumption

Now, its time to put a naive assumption to the Bayes' theorem, which is, **independence** among the features. So now, we split **evidence** into the independent parts.

Now, if any two events A and B are independent, then,

$$P(A,B) = P(A)P(B)$$

Theory Behind Naïve Baiyes

Hence, we reach to the result:

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

which can be expressed as:

$$P(y|x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i|y)}{P(x_1)P(x_2)\dots P(x_n)}$$

Now, as the denominator remains constant for a given input, we can remove that term:

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

Question

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

Now which tag does the sentence
A very close game belong to?

In our case, we have $P(\text{Sports} \mid \text{a very close game})$, so using this theorem we can reverse the conditional probability:

$$P(\text{sports} \mid \text{a very close game}) = \frac{P(\text{a very close game} \mid \text{sports}) \times P(\text{sports})}{P(\text{a very close game})}$$

Since for our classifier we're just trying to find out which tag has a bigger probability, we can discard the divisor—which is the same for both tags—and just compare

$$P(\text{a very close game} \mid \text{Sports}) \times P(\text{Sports})$$

with

$$P(\text{a very close game} \mid \text{Not Sports}) \times P(\text{Not Sports})$$

Question

$$P(a \text{ very close game}) = P(a) \times P(\text{very}) \times P(\text{close}) \times P(\text{game})$$

This assumption is very strong but super useful. It's what makes this model work well with little data or data that may be mislabeled. The next step is just applying this to what we had before:

$$P(a \text{ very close game} | \text{Sports}) = P(a | \text{Sports}) \times P(\text{very} | \text{Sports}) \times P(\text{close} | \text{Sports}) \times P(\text{game} | \text{Sports})$$

Calculating Probabilities

- ❖ $P(\text{Sports})$ is $\frac{3}{5}$. Then, $P(\text{Not Sports})$ is $\frac{2}{5}$
- ❖ Then, calculating $P(\text{game} \mid \text{Sports})$
- ❖ Count of word “game” appears in Sports texts (2) divided by the total number of words in sports (11)
Therefore,
- ❖ $P(\text{game} \mid \text{sports}) = 2/11$
- ❖ “close” doesn’t appear in any *Sports* text!
- ❖ $P(\text{close} \mid \text{Sports}) = 0$

$$P(a|\text{Sports}) \times P(\text{very}|\text{Sports}) \times 0 \times P(\text{game}|\text{Sports})$$

Calculating Probabilities

- ❖ To balance this, we add the number of possible words to the divisor, so the division will never be greater than 1

possible words are ['a', 'great', 'very', 'over', 'it', 'but', 'game', 'election', 'clean', 'close', 'the', 'was', 'forgettable', 'match'].

Calculating Probabilities

$$P(\text{game}|\text{sports}) = \frac{2 + 1}{11 + 14}$$

The full results are:

Word	P (word Sports)	P (word Not Sports)
a	$(2 + 1) \div (11 + 14)$	$(1 + 1) \div (9 + 14)$
very	$(1 + 1) \div (11 + 14)$	$(0 + 1) \div (9 + 14)$
close	$(0 + 1) \div (11 + 14)$	$(1 + 1) \div (9 + 14)$
game	$(2 + 1) \div (11 + 14)$	$(0 + 1) \div (9 + 14)$

Calculating Probabilities

Now we just multiply all the probabilities, and see who is bigger:

$$\begin{aligned} &P(a|Sports) \times P(very|Sports) \times P(close|Sports) \times P(game|Sports) \times \\ &P(Sports) \\ &= 2.76 \times 10^{-5} \\ &= 0.0000276 \end{aligned}$$

$$\begin{aligned} &P(a|Not Sports) \times P(very|Not Sports) \times P(close|Not Sports) \times P(game|Not Sports) \times \\ &P(Not Sports) \\ &= 0.572 \times 10^{-5} \\ &= 0.00000572 \end{aligned}$$

Excellent! Our classifier gives "A very close game" the **Sports** tag.

Part of Speech Tagging

Part Of Speech Tagging

- ❖ Is the task of assigning POS tags to words
- ❖ Selecting among more than one tags that apply
- ❖ Can be used for further NLP tasks
- ❖ Information extraction, Question Answering etc.

Part Of Speech Tagging

- ❖ Since the greeks 8 basic POS have been distinguished
 - ❖ Noun, verb, pronoun, preposition, adverb,
 - ❖ conjunction, adjective, article
- ❖ Modern works use extended list of POS:
 - ❖ 45 in Penn Treebank corpus
 - ❖ 87 in Brown corpus
- ❖ Used for syntactic processing
 - ❖ Speech recognition – Pronunciation may change

Part Of Speech Categories

Closed class. Function words: prepositions, pronouns, determiners, conjunctions, numerals, auxiliary verbs and particles (preposition or adverbs in phrasal verbs)

Open class:

Nouns: people, place and things proper nouns, common nouns, count nouns and mass nouns

Verbs: actions and processes. Main verbs, not auxiliaries

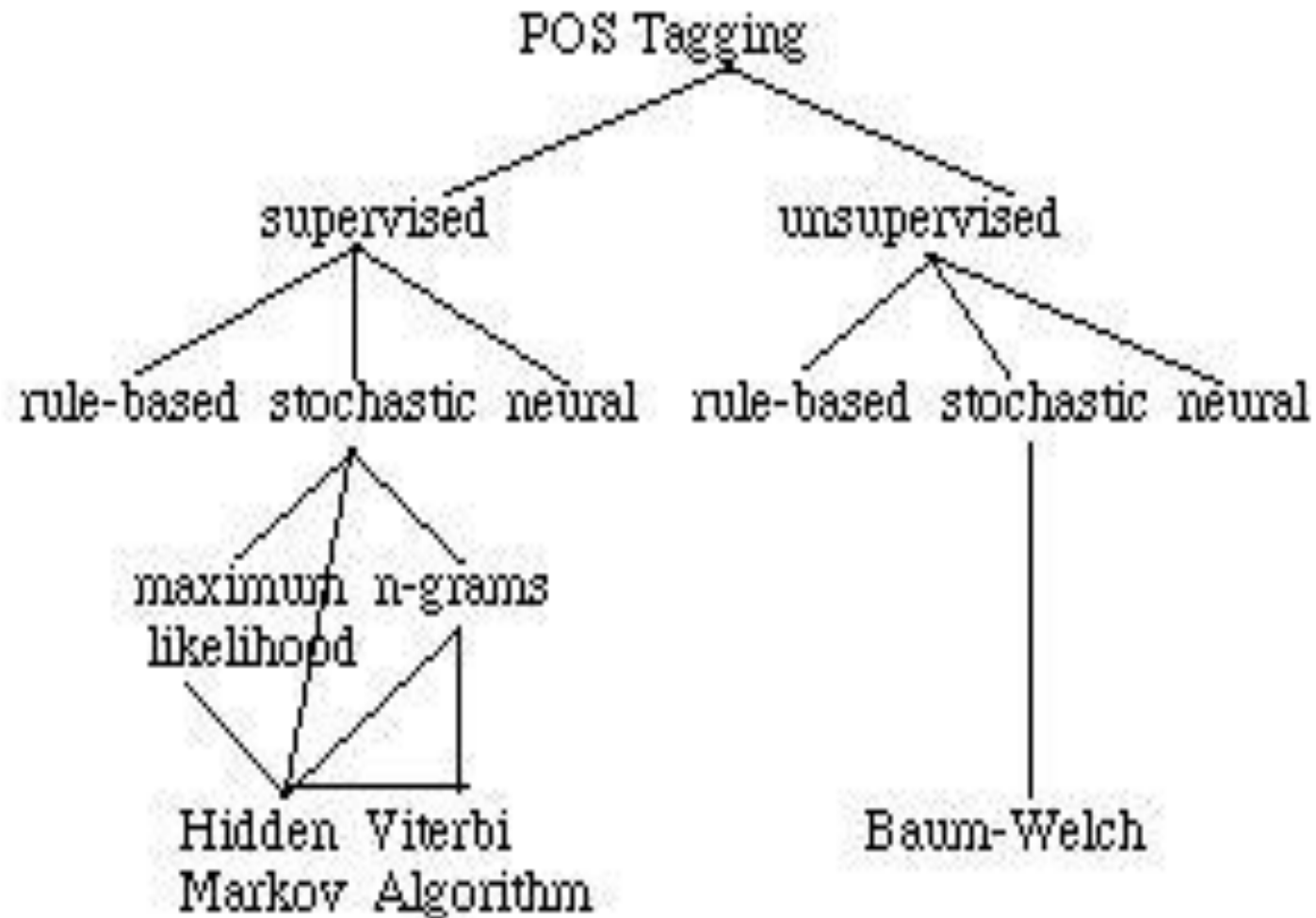
Adjectives: Properties

Adverbs

Existing Taggers

- ❖ Rule based tagging
 - ❖ Brill tagger
- ❖ Stochastic tagging
 - ❖ CLAWS tagger
 - ❖ Tree tagger

Existing Taggers



Rule-Based Tagger

- ▶ The Linguistic Complaint
 - ▶ Where is the linguistic knowledge of a tagger?
 - ▶ Just a massive table of numbers
 - ▶ Aren't there any linguistic insights that could emerge from the data?
 - ▶ Could thus use handcrafted sets of rules to tag input sentences, for example, if input follows a determiner tag it as a noun.

The Brill tagger

56

- ▶ An example of Transformation-Based Learning
 - ▶ Basic idea: do a quick job first (using frequency), then revise it using contextual rules.
 - ▶ **Painting** metaphor from the readings
- ▶ Very popular (freely available, works fairly well)
- ▶ A supervised method: requires a tagged corpus

Brill Tagging: In more detail

- ▶ Start with simple (less accurate) rules...learn better ones from tagged corpus
 - ▶ Tag each word initially with most likely POS
 - ▶ Examine set of **transformations** to see which improves tagging decisions compared to tagged corpus
 - ▶ Re-tag corpus using best transformation
 - ▶ Repeat until, e.g., performance doesn't improve
 - ▶ Result: tagging procedure (ordered list of transformations) which can be applied to new, untagged text

An example

58

- ▶ Examples:
 - ▶ They are expected to **race** tomorrow.
 - ▶ The **race** for outer space.
- ▶ Tagging algorithm:
 1. Tag all uses of “race” as NN (most likely tag in the Brown corpus)
 - ▶ They are expected to **race/NN** tomorrow
 - ▶ the **race/NN** for outer space
 2. Use a transformation rule to replace the tag NN with VB for all uses of “race” preceded by the tag TO:
 - ▶ They are expected to **race/VB** tomorrow
 - ▶ the **race/NN** for outer space

Example Rule Transformations

S	F	B	O		
c	i	o	t		R
o	x	k	h		u
r	e	e	e		l
e	d	n	r		e
-----+-----					
10	18	8	14		NN -> NNP if the tag of words i+1...i+2 is 'NNP'
9	10	1	2		NN -> VB if the tag of the preceding word is 'TO'
8	9	1	18		NN -> VBD if the tag of the following word is 'DT'
7	7	0	9		NN -> VBD if the tag of the preceding word is 'NNS'
6	13	7	8		NN -> JJ if the tag of the preceding word is 'DT', and the tag of the following word is 'NN'
7	9	2	2		NN -> NNP if the tag of the preceding word is 'NN', and the tag of the following word is ','
8	16	8	12		NN -> NNP if the tag of words i+1...i+2 is 'NNP'
4	6	2	11		NN -> IN if the tag of the preceding word is '.'
3	4	1	2		NNP -> NN if the tag of words i-3...i-1 is 'JJ'
3	3	0	2		NN -> JJ if the tag of the following word is 'JJ'
3	3	0	4		NN -> VBP if the tag of the preceding word is 'PRP'
3	3	0	0		WDT -> IN if the tag of the following word is 'DT'

Sample Final Rules

Rules:

```
NN -> NNP if the tag of words i+1...i+2 is 'NNP'
NN -> VB if the tag of the preceding word is 'TO'
NN -> VBD if the tag of the following word is 'DT'
NN -> VBD if the tag of the preceding word is 'NNS'
NN -> JJ if the tag of the preceding word is 'DT', and the tag of the following word is 'NN'
NN -> NNP if the tag of the preceding word is 'NN', and the tag of the following word is ','
NN -> NNP if the tag of words i+1...i+2 is 'NNP'
NN -> IN if the tag of the preceding word is '.'
NNP -> NN if the tag of words i-3...i-1 is 'JJ'
NN -> JJ if the tag of the following word is 'JJ'
NN -> VBP if the tag of the preceding word is 'PRP'
WDT -> IN if the tag of the following word is 'DT'
NN -> JJ if the tag of the preceding word is 'IN', and the tag of the following word is 'NN'
NN -> VBN if the tag of the preceding word is 'VBP'
VBD -> VB if the tag of the preceding word is 'MD'
NN -> JJ if the tag of the preceding word is 'CC', and the tag of the following word is 'NN'
```