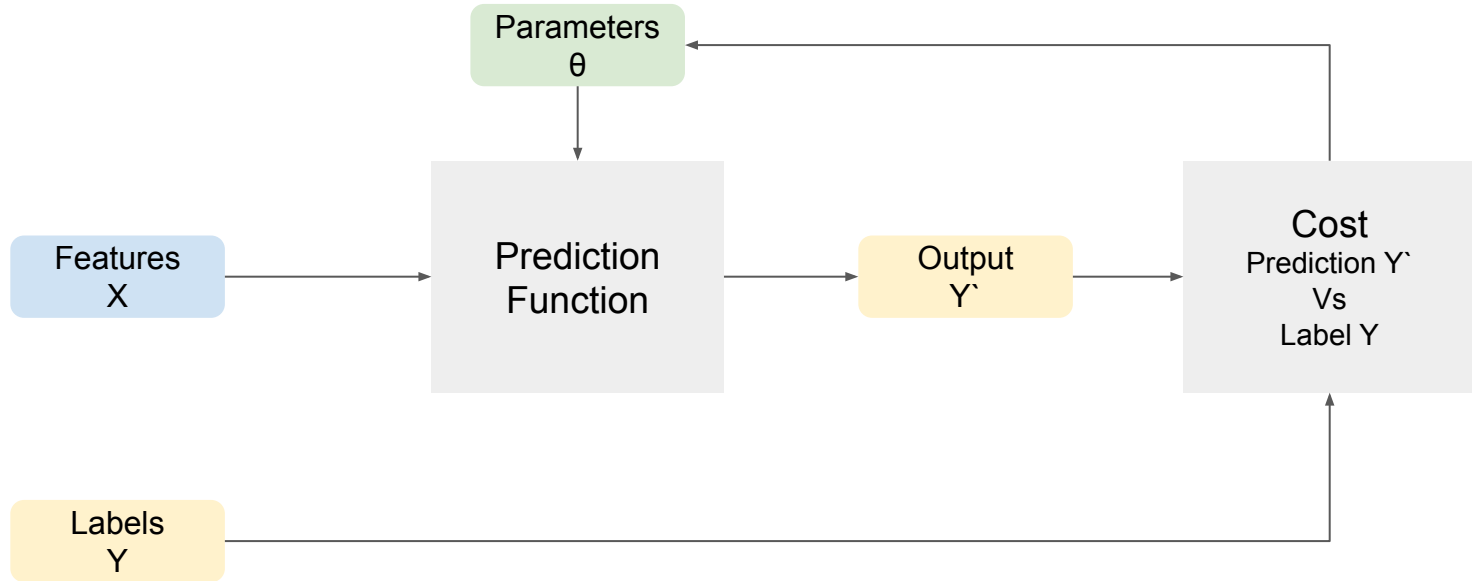


# NLP

Week - 1

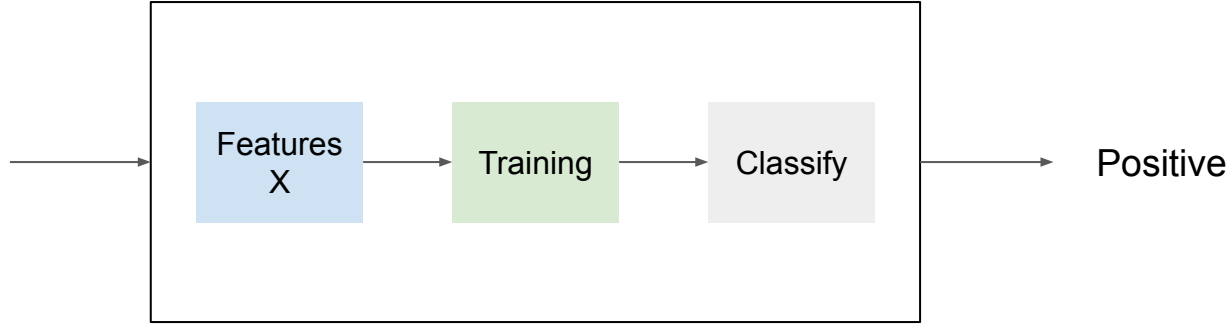
# Supervised Learning



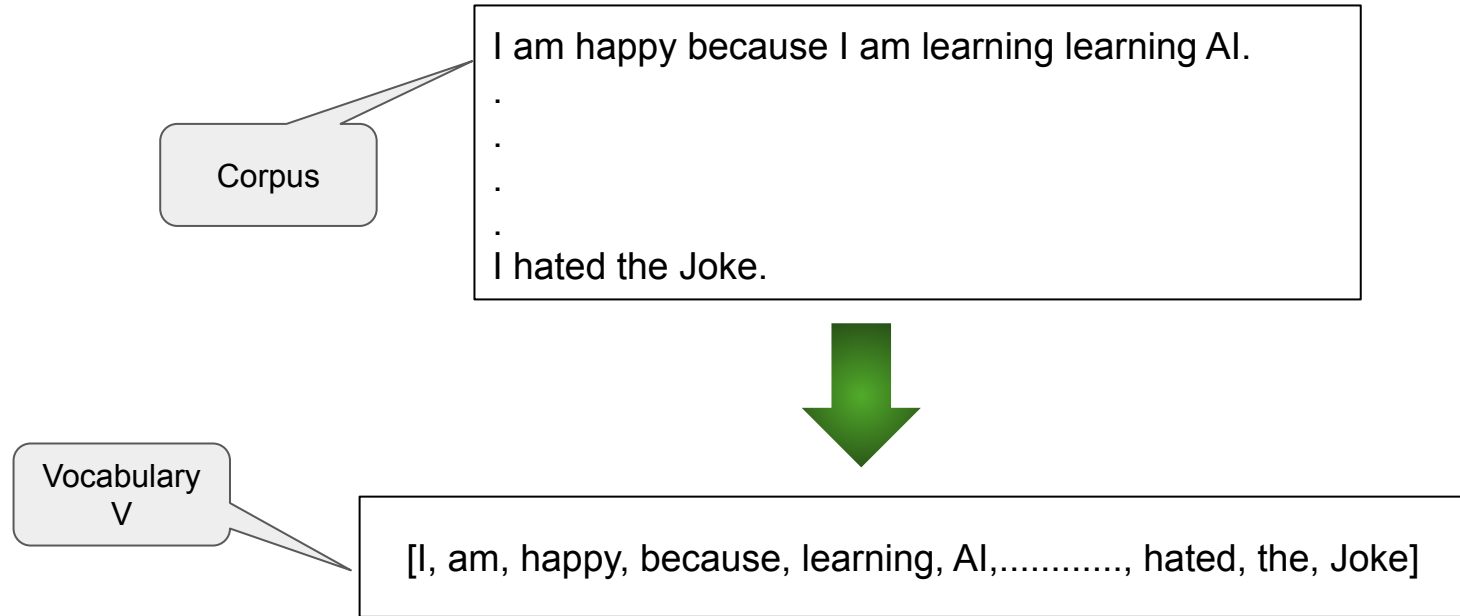
# Sentiment Analysis

I am happy  
because I am  
learning learning  
AI.

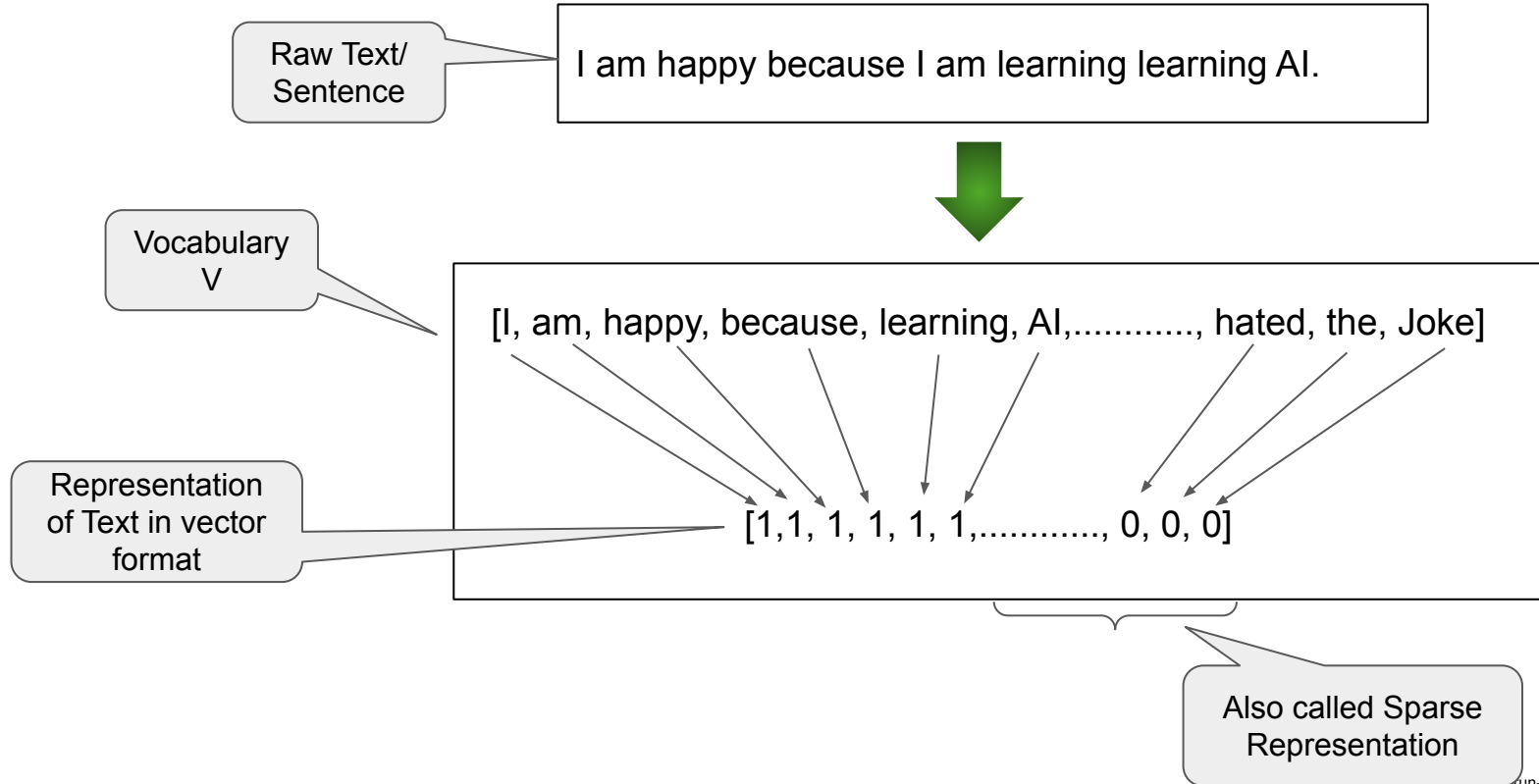
Raw Text



# Vocabulary



# Feature Extraction



# Sparse Representation Limitation

I am happy because I am learning learning AI.



[I, am, happy, because, learning, AI,....., hated, the, Joke]

[1, 1, 1, 1, 1, 1, ....., 0, 0, 0]

$[\theta_0, \theta_1, \theta_2, \theta_3, \theta_4, \theta_5, \dots, \theta_{n-2}, \theta_{n-1}, \theta_n]$

Bias

Too Many Zeros.  
Has to find out  $(Vn + 1)$   
parameters.

- Large Training time
- Large Prediction time

# Another Approach for Feature Extraction

The diagram illustrates the feature extraction formula  $X_m = [ 1, \sum_m \text{freq}(w, 1), \sum_m \text{freq}(w, 0) ]$ . The formula is presented with callouts explaining its components:

- Bias**: A callout pointing to the constant '1' in the first position of the vector.
- Sum of Positive Frequencies**: A callout pointing to the term  $\sum_m \text{freq}(w, 1)$ , which is highlighted in green.
- Sum of Negative Frequencies**: A callout pointing to the term  $\sum_m \text{freq}(w, 0)$ , which is highlighted in red.
- Features of Raw Text m**: A callout pointing to the entire vector  $X_m$ .

# Sentiment Analysis

Corpus

I am happy because I am learning AI.

I am happy.

I am sad, I am not learning AI.

I am sad.

Vocabulary

I

am

happy

because

learning

AI

sad

not



# Sentiment Analysis - Word Frequency in Classes

Positive Corpus

I am happy because I am learning AI.

I am happy.

Vocabulary	Frequency
I	3
am	3
happy	2
because	1
learning	1
AI	1
sad	0
not	0

Arun Kumar Anala

[analaarun.k@gmail.com](mailto:analaarun.k@gmail.com)

<https://www.linkedin.com/in/arun-kumar-anala-35760523/>

# Sentiment Analysis - Word Frequency in Classes

Negative Corpus

I am sad, I am not learning AI.

I am sad.

Vocabulary	Frequency
I	3
am	3
happy	0
because	0
learning	1
AI	1
sad	2
not	1

Arun Kumar Anala

[analaarun.k@gmail.com](mailto:analaarun.k@gmail.com)

<https://www.linkedin.com/in/arun-kumar-anala-35760523/>

# Sentiment Analysis - Word Frequency in Classes

Vocabulary	Positive Frequency	Negative Frequency
I	3	3
am	3	3
happy	2	0
because	1	0
learning	1	1
AI	1	1
sad	0	2
not	0	1

This is called  
**Frequency Dictionary Mapping**  
from (word, class) : Frequency

# Sentiment Analysis - Feature Extraction

Vocabulary	Positive Frequency	Negative Frequency
I	3	3
am	3	3
happy	2	0
because	1	0
learning	1	1
AI	1	1
sad	0	2
not	0	1

Raw Text

I am sad, I am not learning AI.

$$X_m = [ 1, \quad \sum_m \text{freq}(w, 1), \quad \sum_m \text{freq}(w, 0) ]$$

8

11

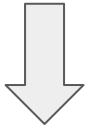
$$X_m = [ 1, \quad 8, \quad 11 ]$$

# Preprocessing : Stop words, Punctuation, Handles, URLs

@Arun @Kabir

Travelling Made Fun. Travel While  
Making Friends.

#ridemate #travelmate #bengaluru  
#delhi #newdelhi



Travelling Made Fun. Travel While  
Making Friends.

Stop  
Words

and  
is  
are  
at  
has  
for  
a  
it

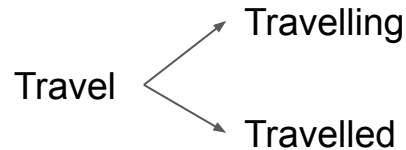
Punctuation

,  
.  
;  
:  
'  
]  
[  
!

# Preprocessing : Stemming and Lowercase

## Stemming

Travelling Made Fun. Travel While  
Making Friends.



## Lowercase

FRIEND

Friend

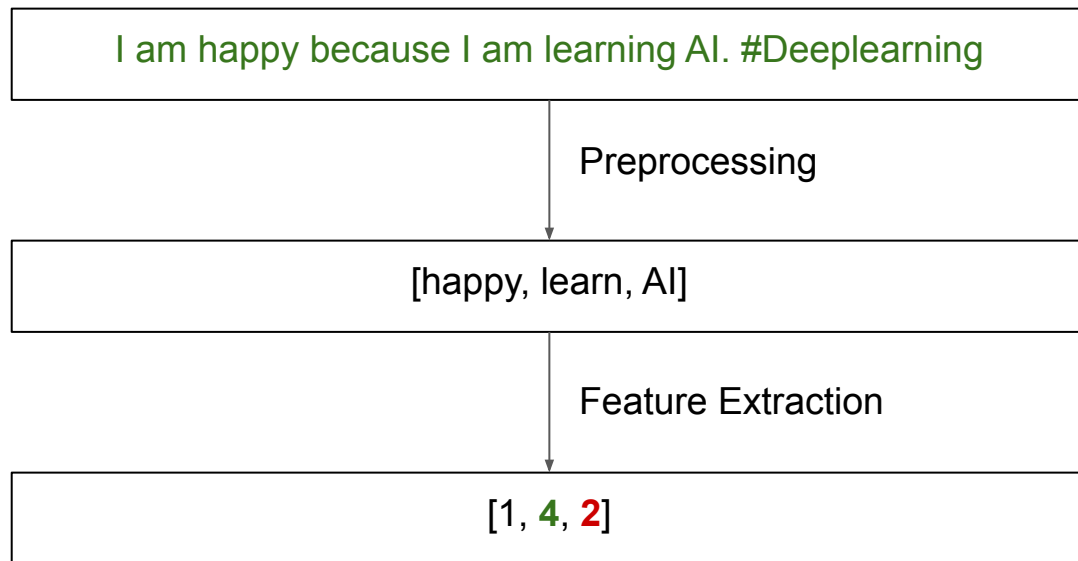
friend

friend

Processed Text

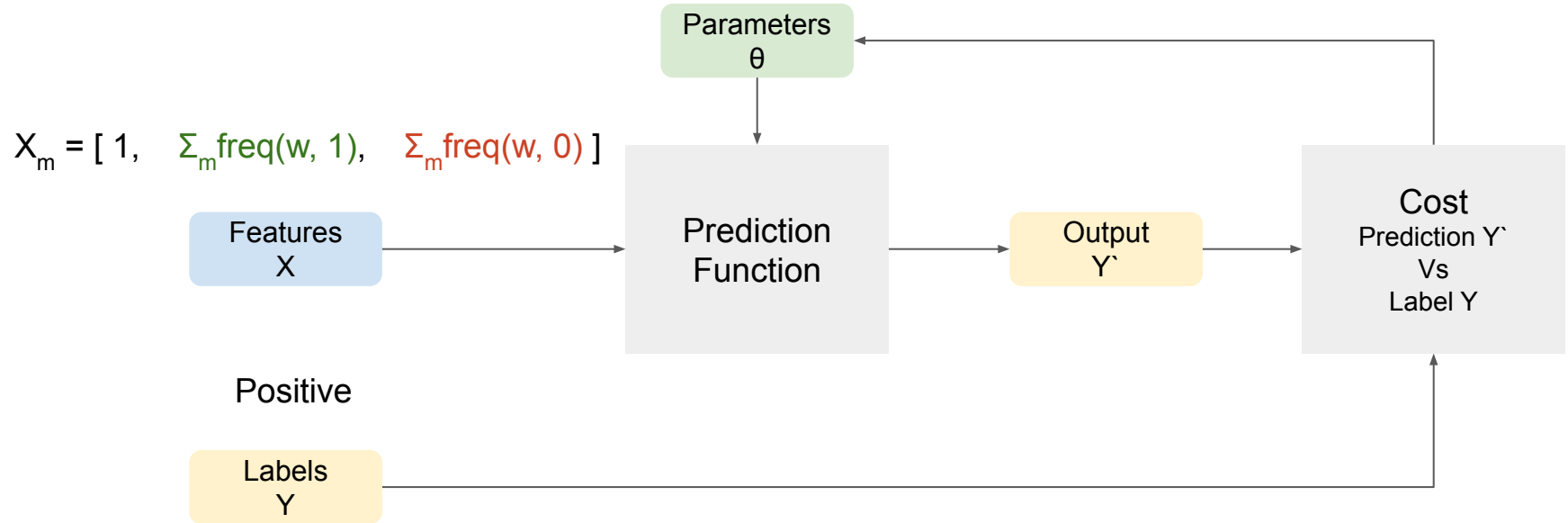
[travel, made, fun, while, making, friend]

# General overview



$$X_m = [ 1, \quad \sum_m \text{freq}(w, 1), \quad \sum_m \text{freq}(w, 0) ]$$

# Supervised Learning





# Probability and Bayes Rule

## Corpus of texts

		Positive		
Negative				

## Texts containing the word happy

	happy	sad
Positive	happy	sad
Negative	happy	sad

# Probability and Bayes Rule

## Corpus of texts

		Positive		
Negative				

$$P(\text{positive}) = N_{\text{pos}} / N = 13/20 = 0.65$$

$$P(\text{negative}) = 1 - 0.65 = 0.35$$

# Probability and Bayes Rule

Corpus of texts

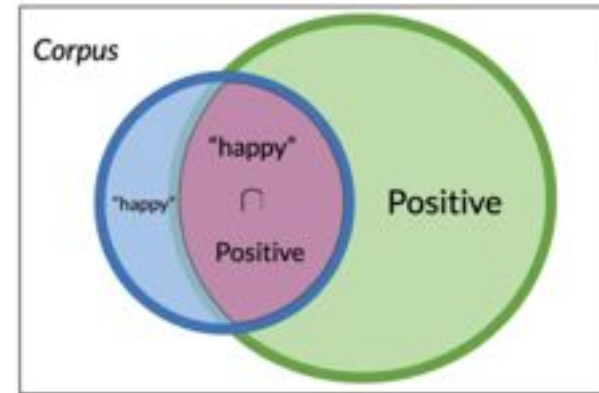
	happy			

$B \rightarrow$  Texts contains “happy”

$$P(B) = P(\text{happy}) = 4 / 20 = 0.2$$

## Probability of Intersection

		Positive	
		happy	

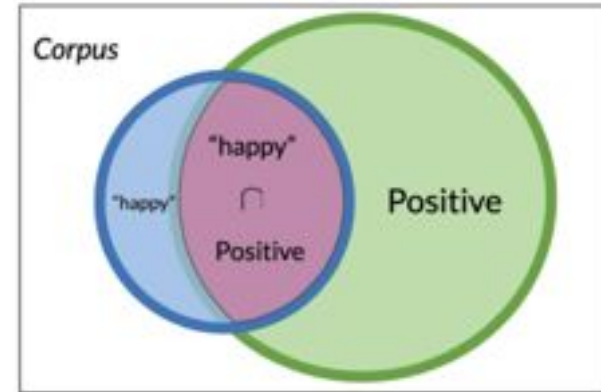


$$P(A \cap B) = P(A, B) = 3 / 20 = 0.15$$

# Conditional Probability

Positive		happy		

$$\begin{aligned} P(A|B) &= P(\text{Positive}|\text{happy}) \\ &= 3 / 4 = 0.75 \end{aligned}$$

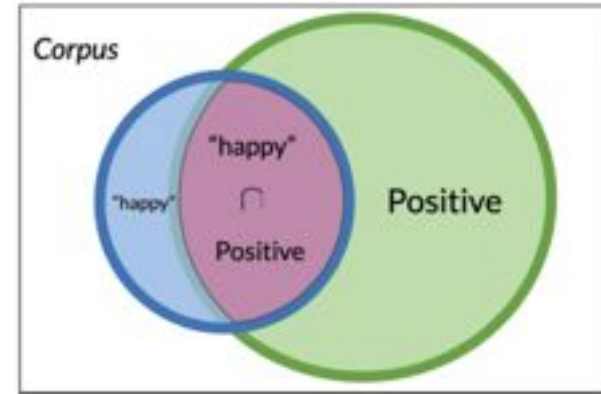


Text has 75% likelihood of being positive,  
if it contains the word "happy"

# Conditional Probability

		Positive		
		happy		

$$\begin{aligned} P(B|A) &= P(\text{happy}|\text{positive}) \\ &= 3 / 13 = 0.231 \end{aligned}$$



Text has 23.1% likelihood of containing the word “happy”, if the text is positive.

# Bayes Rule

$$P(\text{Positive} \mid \text{"happy"}) = P(\text{Positive} \cap \text{"happy"}) / P(\text{"happy"})$$

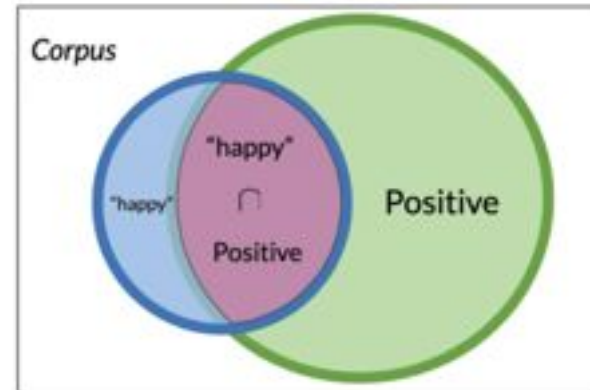
$$P(\text{"happy"} \mid \text{Positive}) = P(\text{Positive} \cap \text{"happy"}) / P(\text{Positive})$$

## Bayes' rule

$$P(\text{Positive}|\text{“happy”}) = P(\text{“happy”}|\text{Positive}) \times \frac{P(\text{Positive})}{P(\text{“happy”})}$$

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

		Positive	
		happy	
		neutral	
		sad	
		Negative	



# Naive Bayes Rule for Sentiment Analysis

Corpus

I am happy because I am learning AI.

I am happy, not sad.

I am sad, I am not learning AI.

I am sad, not happy.

Vocabulary

I

am

happy

because

learning

AI

sad

not

Positive  
Frequency

3

3

2

1

1

1

1

1

13

Negative  
Frequency

3

3

1

0

1

1

2

2

12

Arun Kumar Anala  
[analaarun.k@gmail.com](mailto:analaarun.k@gmail.com)

<https://www.linkedin.com/in/arun-kumar-anala-35760523/>



# Naive Bayes Rule for Sentiment Analysis

Vocabulary	Positive Frequency	Negative Frequency	Positive Probability	Negative Probability
I	3	3	0.24	0.25
am	3	3	0.24	0.25
happy	2	1	0.25	0.08
because	1	0	0.08	0
learning	1	1	0.08	0.08
AI	1	1	0.08	0.08
sad	1	2	0.08	0.17
not	1	2	0.08	0.17
	13	12	1	1

# Naive Bayes Rule for Sentiment Analysis

I am happy today; I am learning

$$\prod_{i=1}^m \frac{P(w_i|pos)}{P(w_i|neg)} = \frac{0.14}{0.10} = 1.4$$

$$\frac{0.20}{0.20} * \frac{0.20}{0.20} * \frac{0.14}{0.10} * \frac{0.20}{0.20} * \frac{0.20}{0.20} * \frac{0.10}{0.10}$$

Vocabulary	Positive Probability	Negative Probability
I	0.24	0.25
am	0.24	0.25
happy	0.25	0.08
because	0.08	0
learning	0.08	0.08
AI	0.08	0.08
sad	0.08	0.17
not	0.08	0.17
	1	1

# Bag of Words

