

# What is sentiment Analysis?

## Reviews

Fragrance-1 (Lavender)	1. smell amazing. 2. must buy super amazing.	Positive
---------------------------	---	----------

## Review

Fragrance-2 (Rose)	1. A decent purchase. 2. Quite okayish.	Neutral
-----------------------	--	---------

## Review

Fragrance-3 (Lily)	1. Smell okay. 2. Decent purchase.	Negative
-----------------------	---------------------------------------	----------

Sentiment Analysis is a use-case of NLP and comes under the category of text-classification.

Sentiment Analysis involves classifying a text

f  
ego dataset,

SNo.	Student feedback	Sentiment label
1.	Timing are very odd such course shouldn't be offered.	negative
2.	Very hard working instructor helps a lot.	positive
3.	Give more programming assign.	neutral

# → How does Sentiment Analysis work?

1. Text - Preprocessing
2. Tokenization
3. Feature Extraction
4. Sentiment Classification

## Types of sentiment Analysis

Document-level

Sentence-level

Aspect-level

→ objective (fact)

→ subjective  
(views of other)

## feature Extraction

• Vector Representation

• Positive & Negative Count

## How does sentiment analysis work?

① Text preprocessing → This step involves cleaning the text data by removing noise, such as punctuation, special characters, and irrelevant information like URLs or HTML tags.

It may also involve tasks like Lowercasing, removing stopwords (commonly used words like "the", "is", "and" etc) and handling word contractions.

② Tokenization → In this step, the preprocessed text is split into smaller units called tokens, which are typically words or phrases. Tokenization makes it easier to analyze and process the text data.

There are different methods of tokenization such as word tokenization which splits the text into individual words, and sentence tokenization which splits the text into sentences.

③ Feature Extraction → Once, the text has been tokenized, features are extracted from the tokens to represent the text data in a format that can be used by ML algorithms.

④ Sentiment Analysis → Finally, ML models or statistical techniques are applied to the extracted features to classify the sentiment of the text. Common approaches includes SVM, logistic, deep learning, RNNs, CNNs.

# NLTK → Natural Language Toolkit

Preprocessing → remove stopwords and punctuations

stopwords      punctuations

and

9

is

0

one

0

at

1

has

4

for

9

2

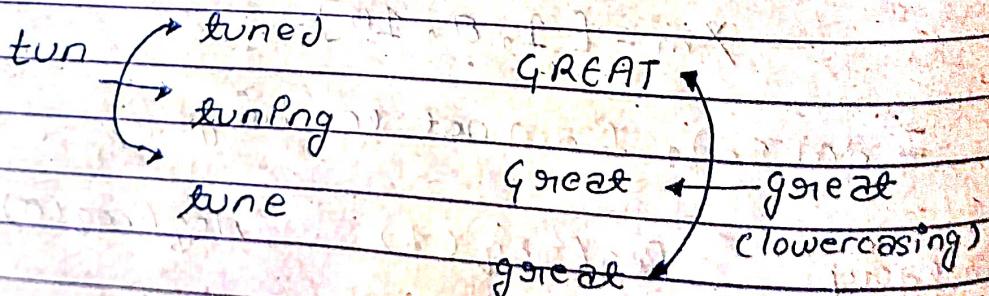
→ @ymouzi and @Andrew\_Ng are tuning

Handles and URL's →

→ @ymouzi and @Andrew\_Ng are tuning

stemming & lowercasing →

tuning great AI model

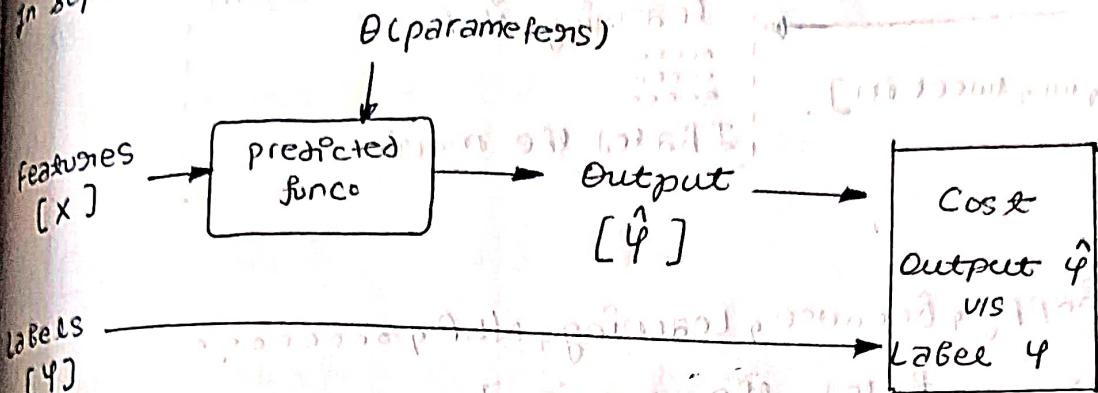


preprocessed tweet ↴

[tun, ai, great, model]

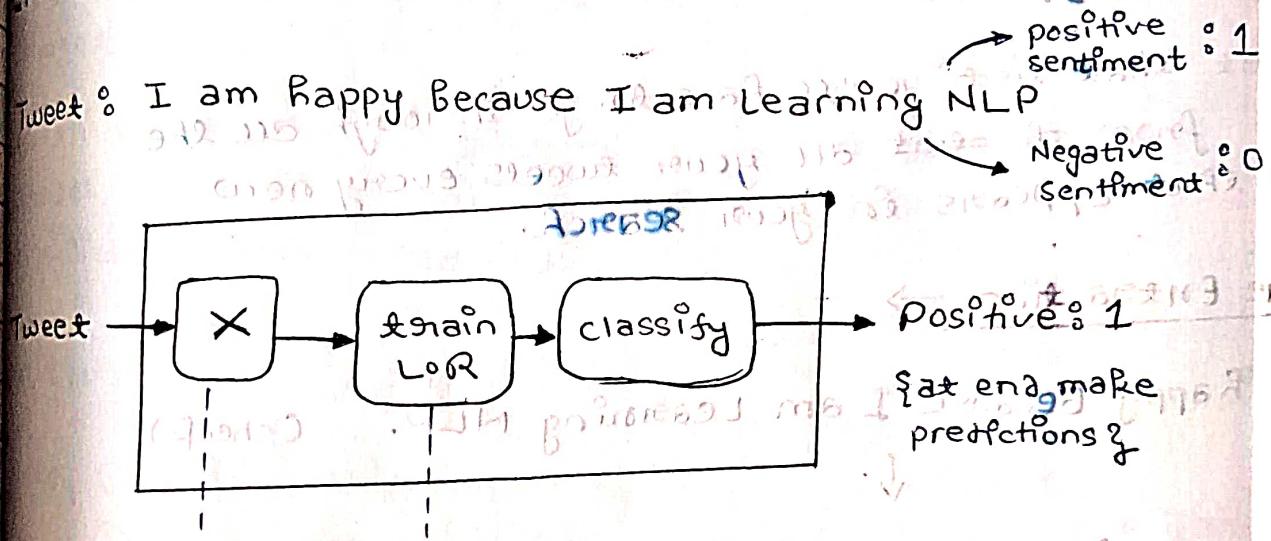
## Logistic Regression for sentiment Analysis of tweets

In supervised ML,



Whether the tweet have a positive or negative sentiment?

Tweet: I am Happy Because I am Learning NLP



Extract useful features while minimizing the costs  
logistic reg classifier

1. Vocabulary & feature extraction:

How to represent a text as a vector →

Build a vocabulary, that will allow you to encode any text or tweet as an array of numbers.

1. Tokenization, 2. Stemming, 3. Stop words removal

1. Tokenization, 2. Stemming, 3. Stop words removal

## Vocabulary

## Tweets

[tweet<sub>1</sub>, tweet<sub>2</sub>, ..., tweet<sub>m</sub>]

I am happy because I am learning NLP.

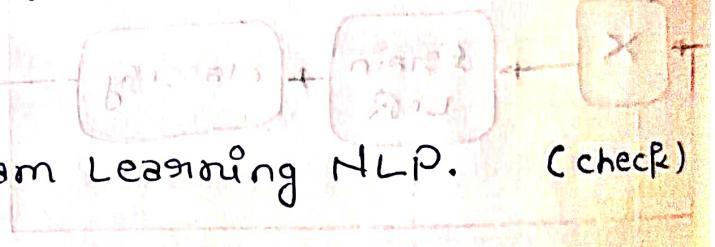
00000  
00000  
00000

I hated the movie.

- Vocabulary (V) is the list of unique words from your list of tweets.
  - To get that list you'll have to go through all the words from ~~the list~~ all your tweets every new word that appears in your search.

### Feature Extraction →

I am happy Because I am Learning NLP. (check)



$V = [I, am, happy, Because, learning, NLP, ---, Rated,$   
 $\text{the, movie}]$ .

- To do so, you will have to check every word from your vocabulary appears for the tweet. You will assign a value to that feature. If it doesn't appear you assign a value



[I am happy because learning NLP, I hated the movie]

$\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$  [1, 1, 1, 1, 1, 1, ..., 0, 0, 0]

6 Count of 1's & many 0's.

systems with "sparse" Representations —

lot of zeros → that's a sparse representation.

I am happy because I am Learning NLP



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

This representation would have a number of features equal to the size of your entire vocabulary.

This would have a lot of features equal to 0 for every tweet.

With the sparse representation, a logistic regression model would have to learn  $(n+1)$  parameters, where  $n = \text{size of your vocabulary}$ .

$$n = |V| \rightarrow [\theta_0, \theta_1, \theta_2, \dots, \theta_n]$$

- ∴  
① Large training time  
② Large prediction time.

### Positive and Negative frequencies —

#### Corpus

I am happy because I am Learning NLP.

I am happy

I am sad, I am not Learning NLP.

I am sad.

#### Vocabulary

I  
am  
happy  
Because  
learning  
NLP  
sad  
not

## Positive tweets

I am happy because I am Learning NLP

I am happy

## Negative tweets

I am sad, I am not learning NLP.

I am sad.

Vocabulary	PosFreq(1)	NegFreq(0)
I	3	3
am	3	3
happy	2	0
because	1	0
learning	1	1
NLP	1	1
sad	0	2
not	0	1

- Previously we encoded a ~~text~~ tweet as a vector of dimension 1 (cv). Now, we'll learn to encode a tweet or specially represented as a vector of dimension 2, 3, ...

freqs : dictionary mapping from (word, class) to frequency

$$X_m = [1, \sum_w \text{freqs}(w, 1), \sum_w \text{freqs}(w, 0)]$$

↓                      ↓                      ↓  
 arbitrary tweet m    sum of pos freqs    sum of neg freqs  
 or  
 features of tweet m

I am sad, I am not learning NLP.



$$\begin{array}{ccccccccc} \text{I} & \text{am} & \text{happy} & \text{because} & \text{Learning} & \text{NLP} & \text{sad} & \text{not} \\ \downarrow & \downarrow \\ 3 & + 3 & + \cancel{x} & + \cancel{x} & + 1 & + 1 & + 0 & + 0 \\ \end{array} \equiv 8$$

$$\begin{array}{ccccccccc} \text{I} & \text{am} & \text{happy} & \text{because} & \text{Learning} & \text{NLP} & \text{sad} & \text{not} \\ \downarrow & \downarrow \\ 3 & + 3 & + \cancel{x} & + \cancel{x} & + 1 & + 1 & + 2 & + 1 \\ \end{array} \equiv 11$$

$$xm = [1, 8, 11].$$

bring it all together →

I am Happy Because I am Learning NLP @deeplearning

↓ [preprocessing]

[happy, learn, nlp].

↓ [feature extraction]

bias ← [1, 4, 2]

↑ ↑

pos freq, neg freq, word responses = 1000 × 100

↓

→ pos freq + neg freq = 1000 × 100 × 2 = [300] %

I am Happy Because I am learning NLP  
@deeplearning

I am sad not learning NLP

.....  
.....

I am sad °(

[happy, learn, nlp]

[1, 40, 20],

[sad, not, learn, nlp]

[1, 20, 50],

[sad]

[1, 5, 35]

$$\begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} \end{bmatrix}$$

[1, 40, 20],

[1, 20, 50],

[1, 5, 35]).

## General Implementation

$\text{freqs} = \text{build\_freqs(tweet, labels)}$

# Build frequencies dictionary

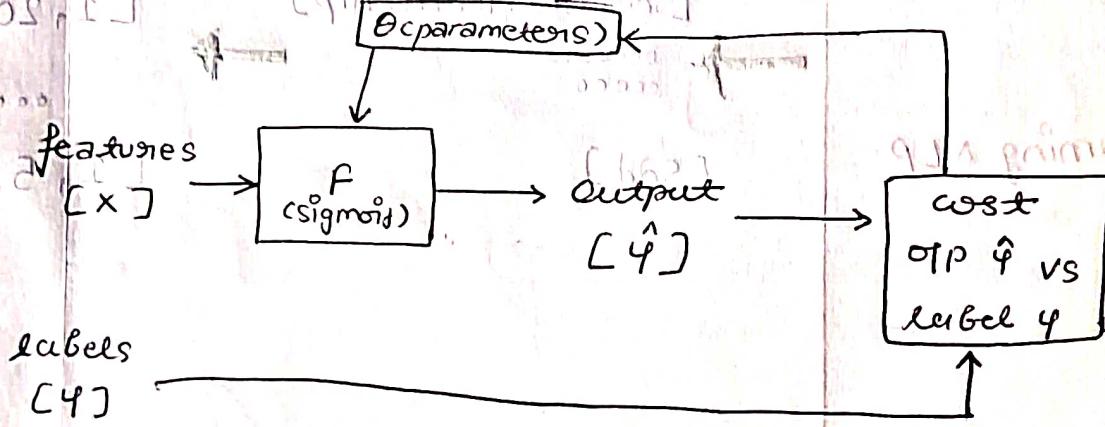
$X = \text{np.zeros}(m, 3)$  # Initialize matrix  $\approx \chi$

for  $i$  in range(m): # for every tweet

$p\_tweet = \text{process\_tweet(tweets}[i])$   
# process tweet

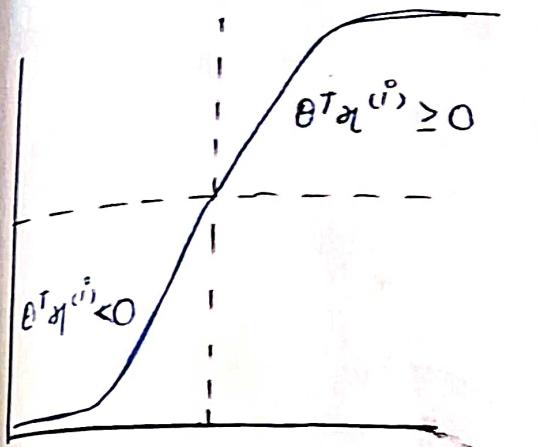
$X[i, :] = \text{extract\_features}(p\_tweet, freqs)$   
# Extract features

## Overview of Logistic Regression →



- To make prediction based on your data, we use a function with some parameters ( $\theta$ ) to map your features to output labels.
- To get an optimum mapping from your features to labels, you minimize the cost function, which works by comparing your output  $\hat{y}$  to the true labels  $y$  from your data.
- After, which the parameters are updated and you repeat this function until your cost is minimized for logistic regression.

$$p_i, e) = \frac{1}{1 + e^{-\theta^T \alpha^{(i)}}}$$



@Yousri and  
@Andrew Ng are tuning a  
GREAT AI model



[tun, ai, great, model]



$$\alpha^{(i)} = \begin{bmatrix} 1 \\ 3476 \\ 245 \end{bmatrix}$$

Sum of all pos & neg freqs of all the  
words in your processed tweets



$$\theta = \begin{bmatrix} 0.00003 \\ 0.00150 \\ -0.00120 \end{bmatrix}$$

