

→ Class start at 9:05 pm

→ NLP (Natural language Processing)

↳ Extract information from text
and take actions on it

→ Convert text into numerical
value and convey all the meaning
in sentence

→ Understand the meaning of text

6		
	255	119
1		

Ex: 1) The food is very good ✓

Sentiment
+ve

✓ 2) Completely lacking in good taste
good service

-ve

Applications of NLP

- Translation :
- Chatbot
- Generate a summarized review
- Sentiment analysis
- Q & A

Syllabus

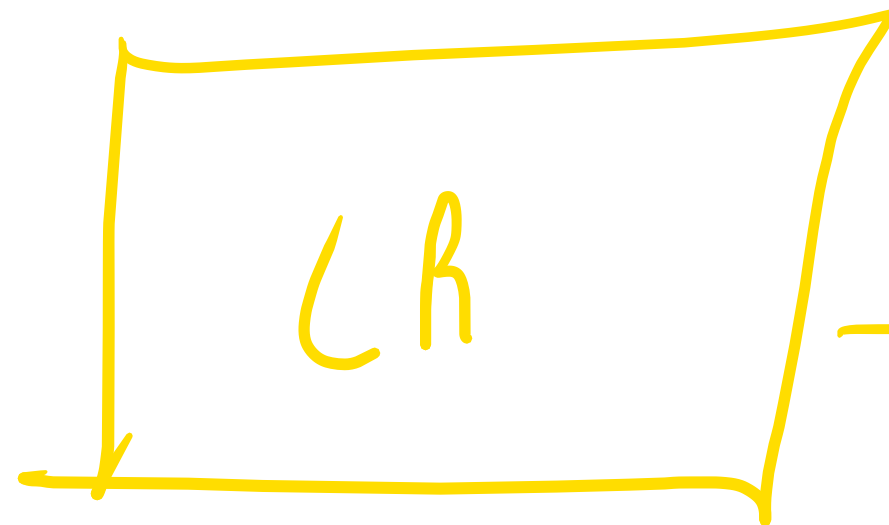
- Intro to NLP
 - Text Representation
- } Preprocessing, Encoding, Bow, TF-IDF
- Word Embedding → CBOW, Skip Gram
 - Language Modelling
 - RNN, LSTM, GRU
 - NER
 - Transformer, Attention, BERT

Problem

→ Given the corona - tweet,
classify the sentiment.

text

X - features



Sentiment %

Preprocessing technique

→ Removal of special characters

Sentiment → ~~https://~~ _____ .com
😊 →

Ticket → IT
→ Engineering 😊

→ Removal of stopword

↳ Commonly used word in any language
↳ a, an, the, and

Corpus : A large & structured collection of text document. Can be book, article, collection of social media post.

Document : Collection of large document is called corpus



Vocabulary

↳ A set of unique words in corpus

1. The movie is amazing. I watched it today

2. I watched the movie Spiderman

Vocabulary = { the, movie, is, amazing, I, watched
it, Spiderman }

Tokenization

↳ Breaking text into smaller units called tokens

✓ ~~↳ Word~~ →

↳ Character

✓ ~~↳ Subword~~

[I, the, movie, is, amazing,
I, watched, it, today]

→ Break until 10:10pm

Stemming & Lemmatization

→ Helps in converting words into their base form

Eg: have, having, had → have
run, running, ran → run

Stemming

history — histor
historical — histor

→ May not be understood by human

→ Works by cutting begin or end of word

Lemmatization

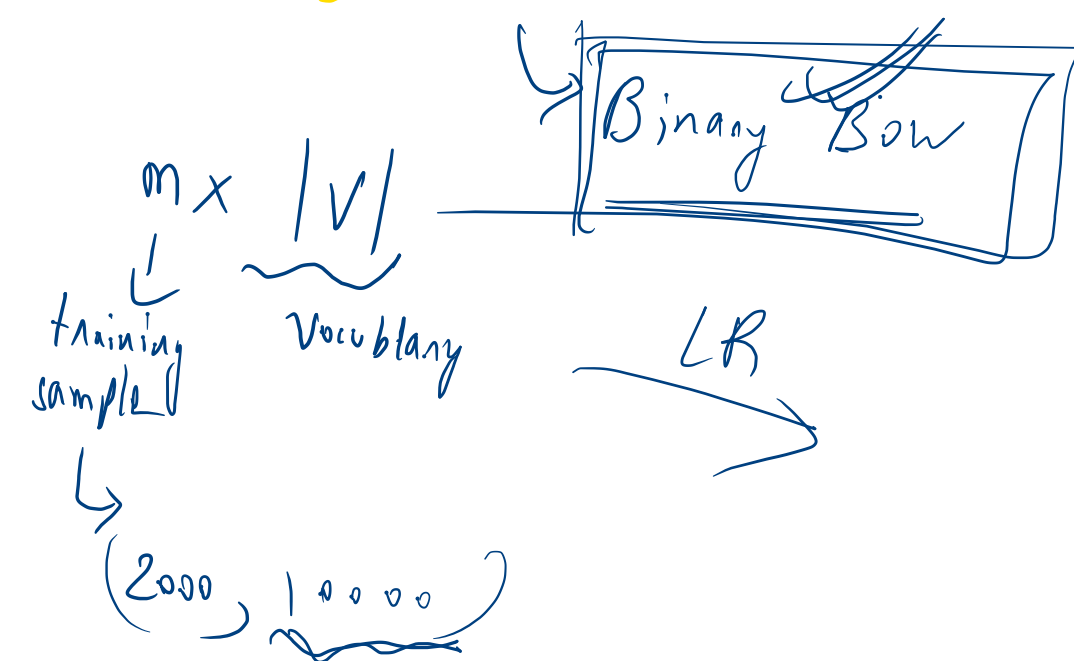
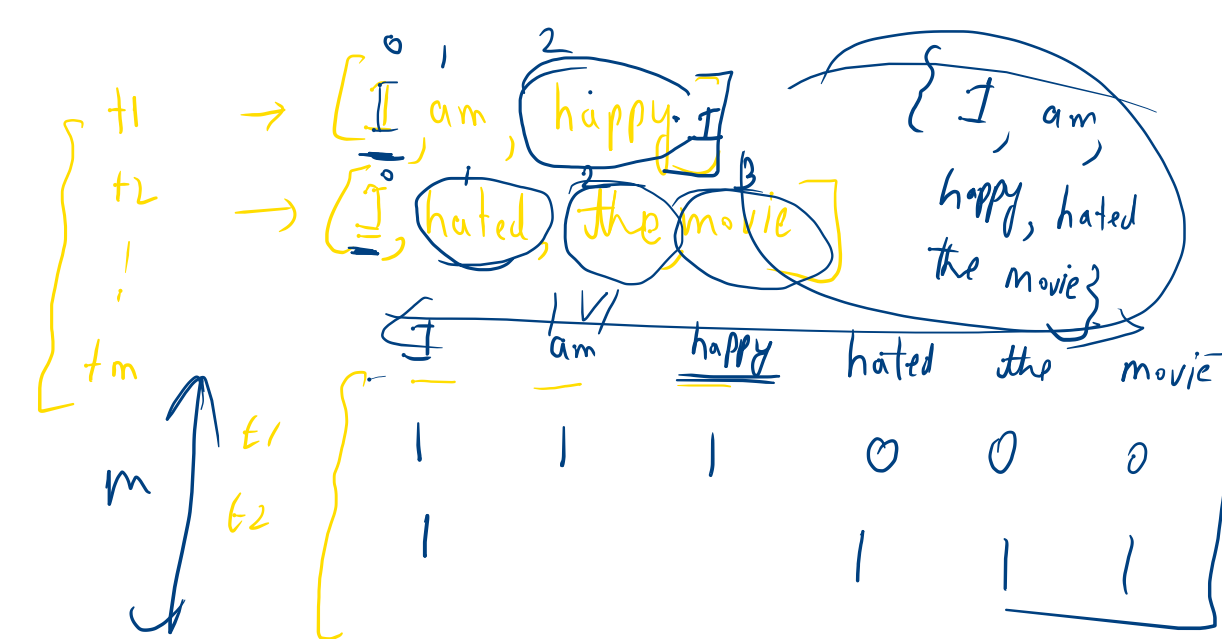
history > history
historical > history

→ Understood by human

→ Considers the context and convert into more meaning full

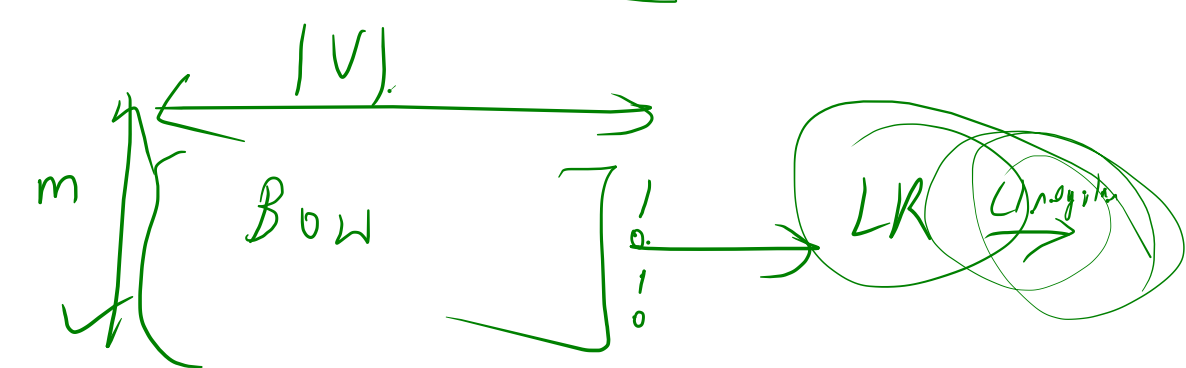
1. Preprocess
2. Tokenize
3. Stem & Lemm

4. [- - -]



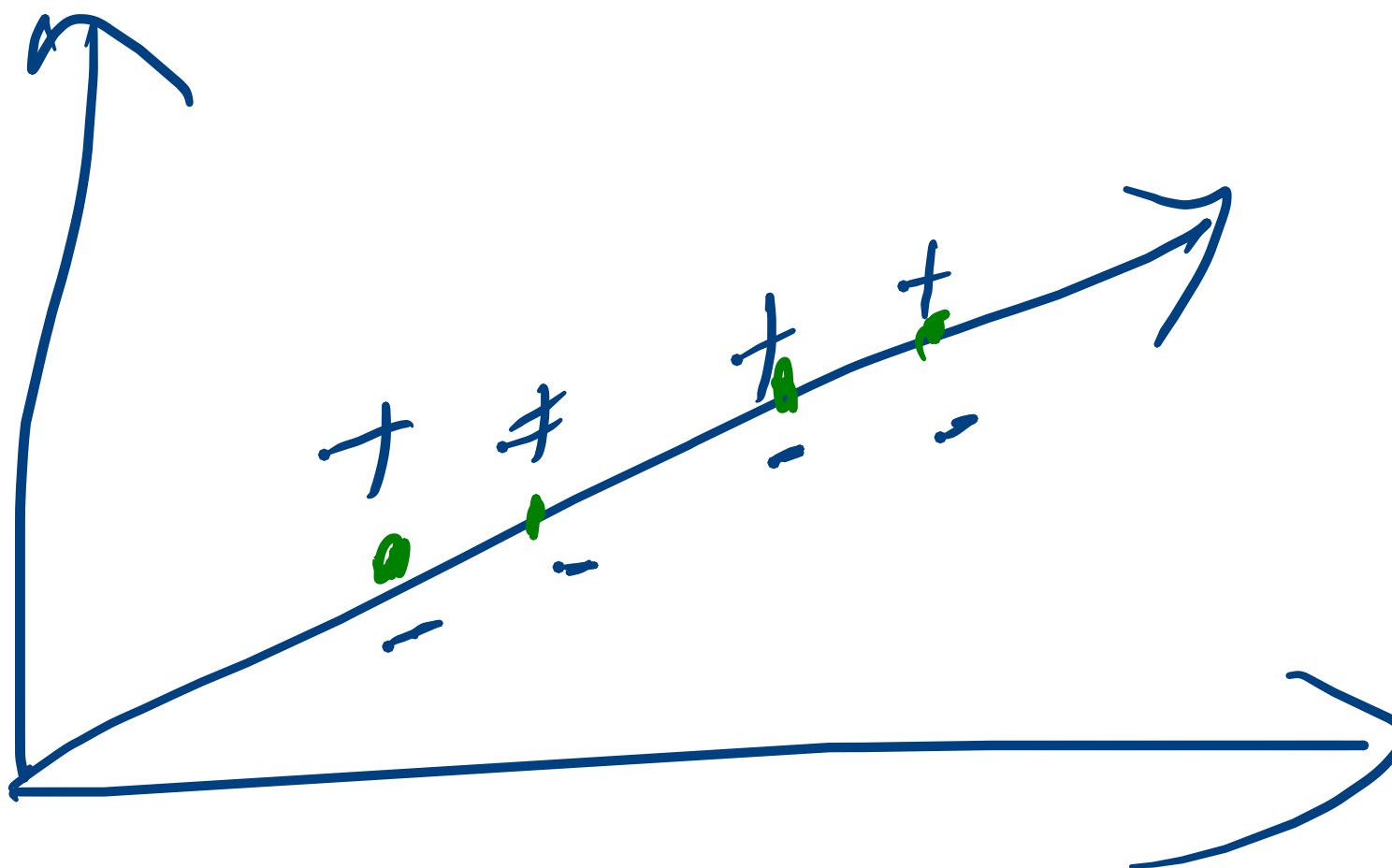
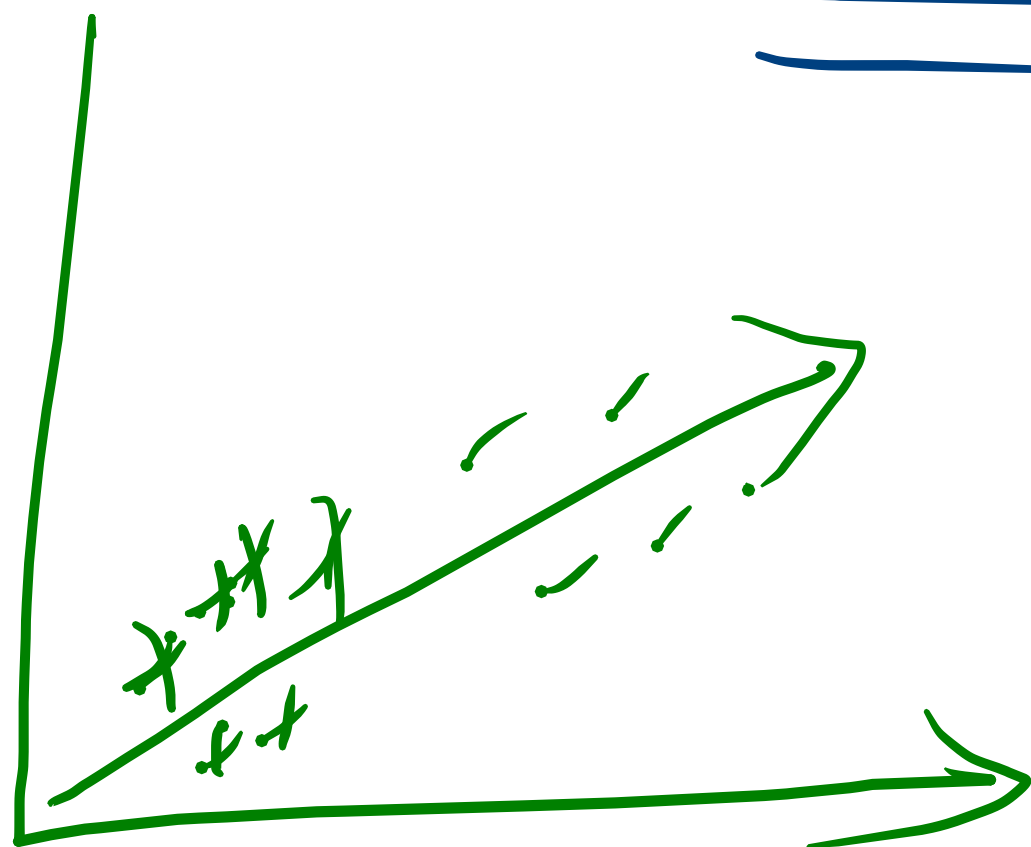
→ Very sparse

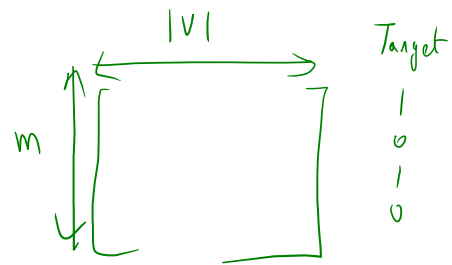
→ LR will have $(|V|+1)$ parameters



→ How can we reduce dimension?

↳ PCA



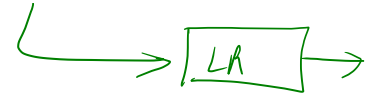


freqs: dictionary mapping from (word, class) to frequency

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

Features of tweet m Bias Sum Pos. Frequencies Sum Neg. Frequencies

2 features



1. I am happy because I am learning NLP
 2. I am happy
 3. I am sad, I am not learning NLP
 4. I am sad
- Target
1
0
0
0

Vocab = { I, am, happy, because, learning, NLP, sad, not }

$$(I, 1) = 3 \quad (am, 1) = 3$$

$$(I, 0) = 3 \quad (am, 0) = 3$$

Surprise

I am sad, I am not learning NLP

(1, 2)

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

$$\sum_w \text{freq}(w, 1) = \text{freq}(I, 1) + \text{freq}(am, 1) + \text{freq}(sad, 1) + \text{freq}(not, 1) + \text{freq}(learn, 1) + \text{freq}(NLP, 1)$$

$$= 3 + 3 + 0 + 0 + 1 + 1 = 8$$

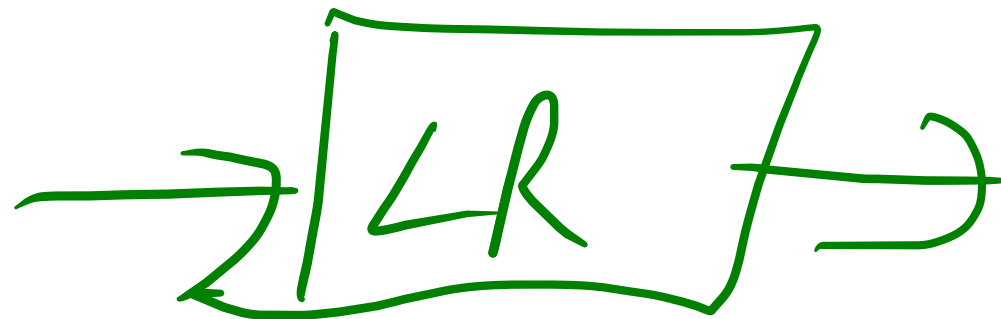
$$\sum_{\omega} f_{req}(\omega, 0)$$



$$m \times 1/1/1$$

$$\sum f_{req}(u)$$

$$\sum f_{req}(u, j) \quad m \times 2$$



m

