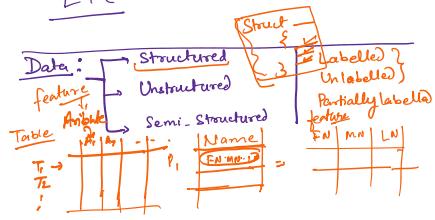


- Problems:
- Classification → Spam filter, sentiment Analysis
  - Regression
  - Ranking
  - Machine Translation (Source Language → Target L)
  - Text Summarization
  - NER
  - Ent. rec., Name, Organi/Place etc.
  - Part-of-speech Tags in (POS) Who, What, Why, How
  - Q&A
  - ImageCaptioning: Image → Text
  - Image Summarization
  - Text to Image → Text → Image

- Text to speech Recom'
- Speech to Text Translat.
- Speech Data / Audio data → Methods of Text you can apply
- Video Summarization → Embedding Audio to Videos → Seg. Obj. → Seg. Obj. → Summary in Text → Video

ETC.Data

Sample  $\langle f_1^1, f_2^1, \dots, f_d^1 \rangle$  → rows → data points | samples

Sample  $\langle f_1^n, f_2^n, \dots, f_d^n \rangle$

2-features  
 $S_1 \rightarrow \langle f_1^1, f_2^1 \rangle$   
 $S_2 \rightarrow \langle f_1^2, f_2^2 \rangle$   
 $S_3 \rightarrow \langle f_1^3, f_2^3 \rangle$

Table 3-features  
 $S_1 \rightarrow \langle f_1^1, f_2^1, f_3^1 \rangle$   
 $S_2 \rightarrow \langle f_1^2, f_2^2, f_3^2 \rangle$   
 $S_3 \rightarrow \langle f_1^3, f_2^3, f_3^3 \rangle$   
 $S_4 \rightarrow \langle f_1^4, f_2^4, f_3^4 \rangle$   
 (2-dim preprocessed/database view)

→ Unstructured Data 96%

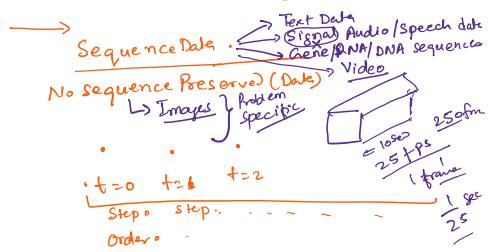
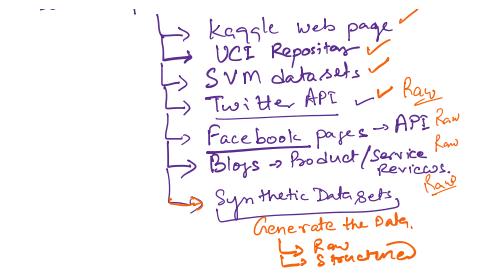
↳ Text, Audio, Image, Video  
 → Don't have the track of feature locations  
 $S_1 \rightarrow$  Mohan is reading a book  
 $S_2 \rightarrow$  Reading a book is very hard

→ Semi-structured → Structured  
 Supervised | Ch. Supervised | Semi-Supervised  
 Labelled | Unlabelled | Partially Labelled  
 $\langle S_i, Y_i \rangle$  |  $\langle S_i \rangle$  |  $\langle S_i \rangle$   
 Problem | Inter. Prob. | Tagging  
 $\langle S_i, f_i \rangle$  |  $f_i \rightarrow$   $\langle S_i \rangle$  |  $\langle S_i : \text{Summary} \rangle$   
 $\langle S_i \rangle$  |  $\langle f_i : \text{Transl. in Target} \rangle$

Assignment: find out the publicly available data in all the categories.

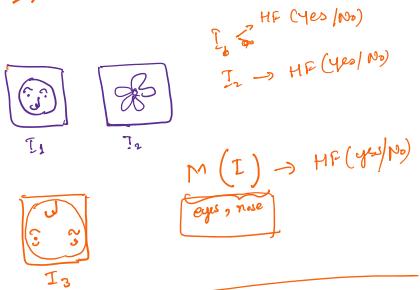
20/1/2021 → Source for the datasets:

- ↳ Kaggle Web page ✓
- ↳ UCI Repository ✓
- ↳ SVM datasets ✓
- ↳ TREC-APT ✓ Ray



HAPPY

H A P P Y      I. am teaching NLP Course.



22/01/2021 Lab session

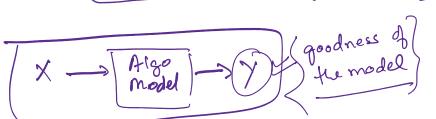
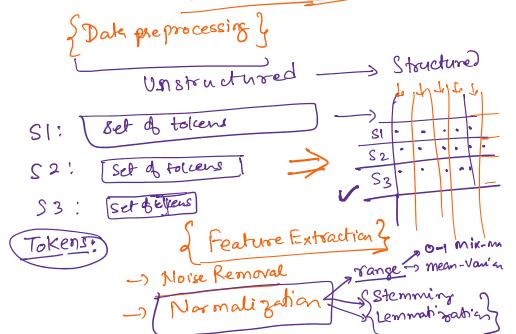
→ Get the Data set ↳ Growl (API)  
 ↳ Download the News Articles  
 ↳ Audio (Speech)

Unstructured Data ↳ Raw form

→ NLP\_Lab Folder → LabAssignments  
 ↳ But raw Data → Raw Data

Browser / Write some script in any of  
 favourite programming lang.

Unstructured ↳ Transformations  
 Structured forms

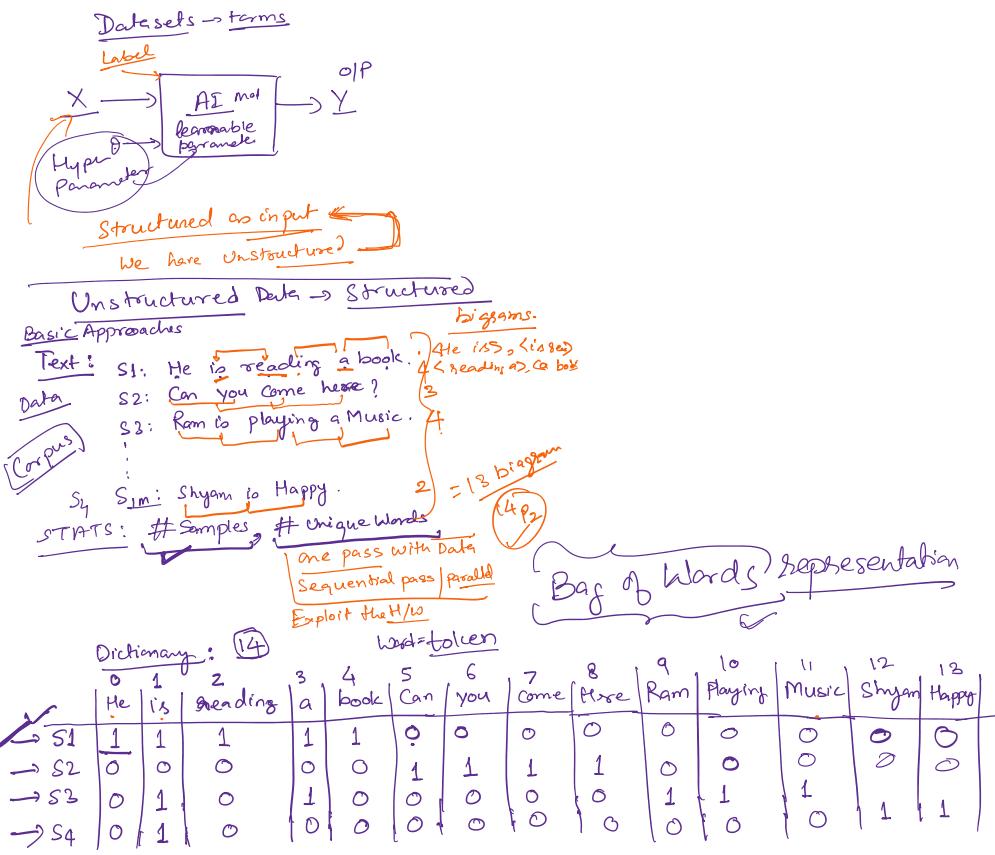


[Group Name: Members]  
 Lab ↳ Get the data ↳ set problems w/ this data

Part (a) ↳ Part (b) ↳ Data cleaning + Data Preprocessing

25/01/2021:

Datasets → forms  
 Label ↳ AI mod ↳ O/P



Token: Word = Feature/dim

Lost: Context ordering

S1:  $\langle (0, 1), (1, 2), (2, 3) \dots \rangle$   
S2:  $\langle \dots \rangle$

Uni-gram, bi-gram model,

2 words:  $\langle (0, w_1), (1, w_2) \rangle$

$\langle w_1, w_2 \rangle, \langle w_2, w_1 \rangle$

$\langle w_1, w_2 \rangle, \langle w_2, w_3 \rangle, \langle w_1, w_3 \rangle, \langle w_2, w_2 \rangle$   
 $\langle w_2, w_1 \rangle$

n-words:  $\binom{n}{2} \Rightarrow n(n-1) / 2$

$$\begin{aligned} & n + n(n-1) / 2 = O(n^2) \\ & \approx 1m + 1m \times 1m = O(m^2) \\ & \binom{n}{1} + \binom{n}{2} + \binom{n}{3} + \dots + \binom{n}{k} = O(n^k) \end{aligned}$$

Theory of ASSIAN ENT

Worst Case Size of Dictionary

Large No. of columns in the data  
High dimensions & Large # features

Highly sparse Data

Word Stem

Prefixes, Suffixes, rules

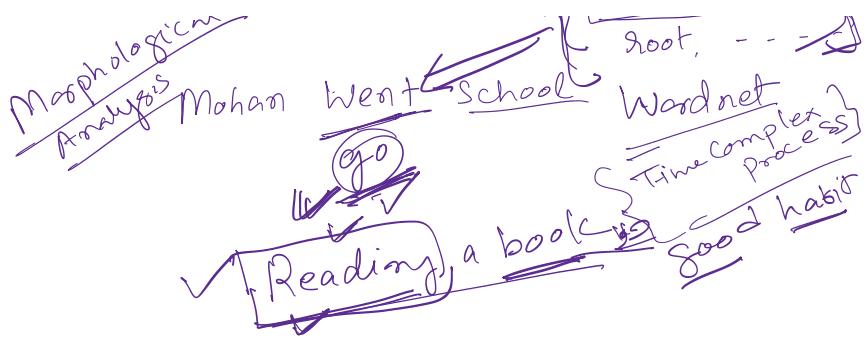
Stemming (Rule Based)

Lemmatization

Database

Root, - - -

Words → base form  
Morphological  
no Mahan Went School Inland net



✓ Root + POS  
Syntactic Enhanced Meaning

Went (Verb, Past)  
 $\Rightarrow$  (go) + Verb + Past

✓ Prex + Base + Suffix

Passes ✓

Lab Assignment:  $X = \{x_i\}_{i=1}^n$

- - 56 - - -  
 five Six  
 number

→ 28 | 01 | 2021      Unstructured

Binary of  
 Text Representation in a Structured Way

	<sup>Term 1</sup> W <sub>1</sub>	<sup>Term 2</sup> W <sub>2</sub>	<sup>Term 3</sup> W <sub>3</sub>	<sup>Term 4</sup> W <sub>4</sub>	<sup>Term 5</sup> W <sub>5</sub>	...	Term m
S <sub>1</sub>	0	4	0	2	1	0...	1
S <sub>2</sub>	1					...	
S <sub>m</sub>	1					...	

1. Binary Representation → Vector of 0's & 1's
2. Term Frequency Rep. → Vector of Integer including zeros.

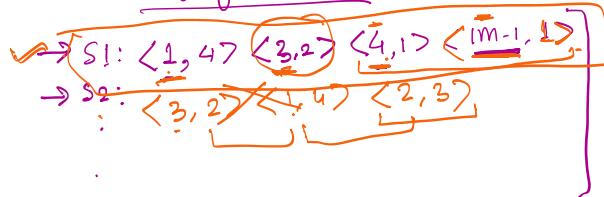
↳ Uni-gram {each word is a term?}  
 ↳ Bi-gram {each pair of words is a token?}

↳ K-gram {group K-words will form a token?}

Highly sparse

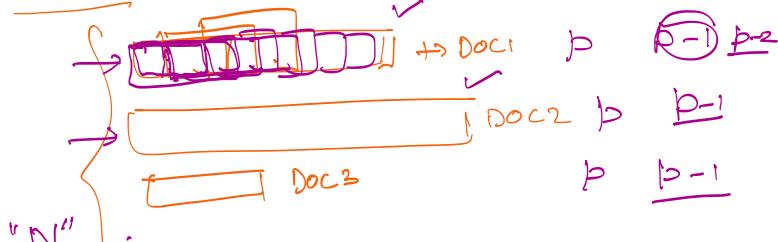
↳ K-gram > group K-words into

Highly sparse



SVM format:  $\langle fid: value \rangle \langle fid_{i+1} \dots \dots \rangle$

3-gram



A average length of a document is 'p'. words

Dictionary size w̄ Uni-gram is 'n'

# tokens = 'n'

# Docs = 'N'

Unigram model

Expected # token =  $Np$  if each

sentence has unique words

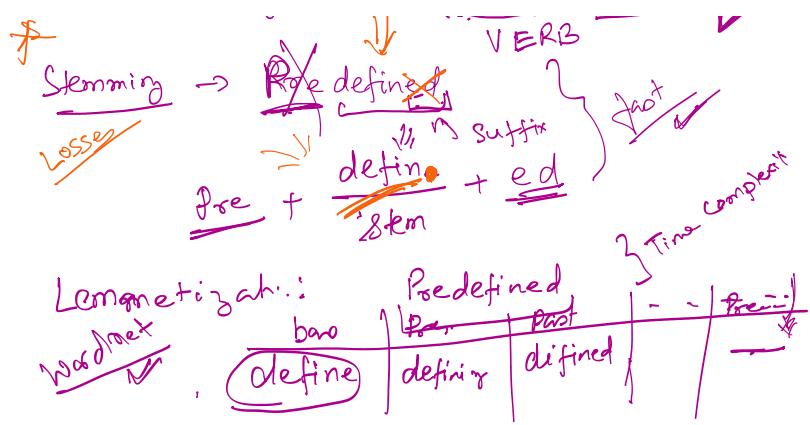
$$NP = n$$

$$n \leq NP$$

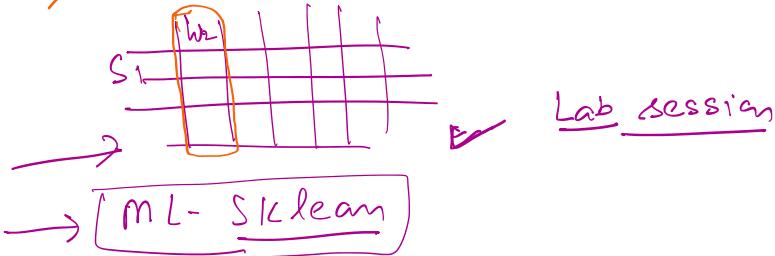
$$\left\{ \begin{array}{l} n_b \leq N(p-1) \\ n_t \leq N(p-2) \\ n_k \leq N(p-k+1) \end{array} \right.$$

→ Raw words/terms  
→ Preprocess it →  
Reading → Read  
Went → go  
Playing → Play

Affixes, base word / root word + POS  
Reading → ~~ds~~ Read VERB  
Present Continuous form  
Stemming → ~~Role defined~~ Root

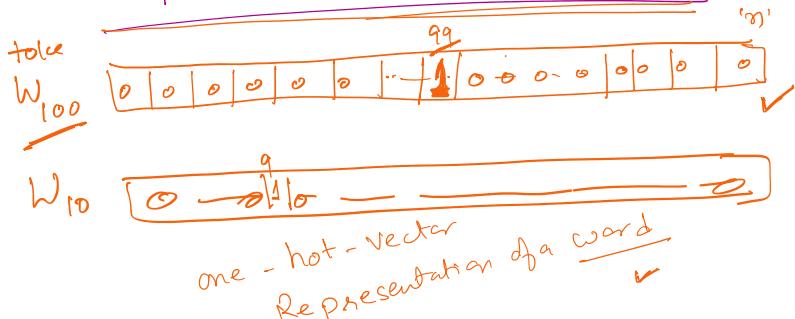


- Read Data
- Tokenize the Data
- Remove Stop Words
- Stemming | Lemmatization
- Vector Representation ↗ Dictionary ↗ Create the Data



Dictionary of 'n' tokens

→ Represent a token in vector form



{ Similar Words } How will you do that?  
 { Group Words } ???

S1..  $w_1, w_2, w_3, w_4$  ↗ Sum of  
 S1:  Term frequency

{ 4 \* 0 1 0 | ... | 0 }  
 ... 0 0 0 1 0 0 0 0

{	$\begin{matrix} 0 &   & 1 &   & 0 &   \\ \downarrow & & \downarrow & & \downarrow & \\ 2x & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{matrix}$
	$\vdash \begin{matrix} 0 &   & 0 &   & 0 &   & 1 &   & 0 & \dots &   & 0 \end{matrix}$
	$\vdash \begin{matrix} 0 &   & 0 & 0 & \dots & \dots &   & 1 \end{matrix}$

~~Term Frequency~~  $S1: \boxed{0|4|0|2|1|0|0\dots|0|1}$  } ~~Term Inc~~

$S1:$   $\boxed{0|4|-\dots-|0|}$   $\stackrel{1m}{\overbrace{\quad\quad\quad}}$   $\stackrel{1m}{\overbrace{\quad\quad\quad}} - -$

$S1: I \text{ am happy } \stackrel{1m}{\overbrace{I\quad\quad\quad}} \stackrel{1m}{\overbrace{am\quad\quad}} \stackrel{1m}{\overbrace{happy}}$   $\stackrel{3m}{\overbrace{\quad\quad\quad}}$

~~Binary~~ ✓  $S1: \boxed{I} \stackrel{1m}{\overbrace{\quad\quad\quad}} \boxed{am} \stackrel{one-m}{\overbrace{\quad\quad\quad}} \stackrel{1m}{\overbrace{happy}}$

$S2: \text{Mohan is reading NLP Books}$

✓  $S2: \boxed{\text{Mohan}} \stackrel{1m}{\overbrace{\quad\quad\quad}} \boxed{\text{is}} \stackrel{1m}{\overbrace{\quad\quad\quad}} \boxed{\text{reading}} \stackrel{2m}{\overbrace{\quad\quad\quad}} \boxed{\text{NLP}} \stackrel{2m}{\overbrace{\quad\quad\quad}} \boxed{\text{Books}} \stackrel{3m}{\overbrace{\quad\quad\quad}} \stackrel{5m}{\overbrace{\quad\quad\quad}}$

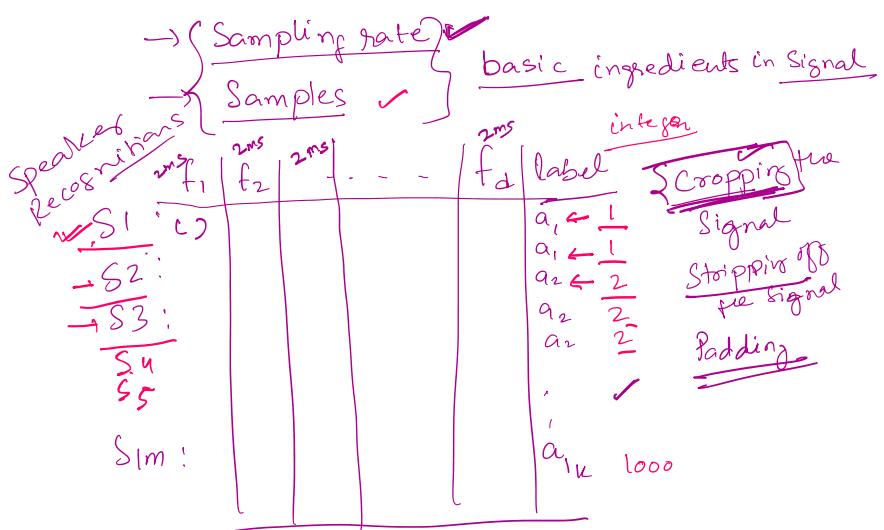
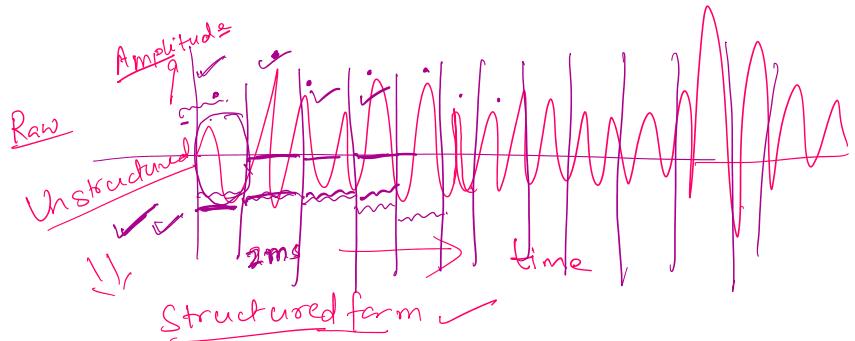
→ ~~# tokens > K~~ → Discards  
~~# tokens < K~~ → Pad it with Zero Vectors

$S1: \dots \boxed{\text{Predefined}} \dots$  ↓  
✓ ~~define~~  $\Rightarrow$  Predefined, defined  
defining, predefined  
Explain!! ↑

Feb 01, 2021		Raw Data (Text)
↓		Term Frequency
$w_1$	$w_2$	$\dots$
$w_d$		label
$S1$		+1
$S2$		+1
⋮		-1
⋮		-1
		⋮

$\left. \begin{array}{l} \text{Train 2-PT} \\ \text{Test} \end{array} \right\}$

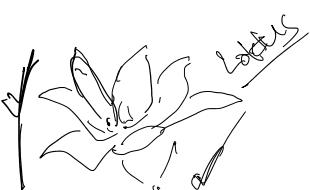
## Audio / speech / Signal



$$\text{Model}(S_{\text{test}}) \rightarrow \underline{596}$$

## Image Data:

8-bit 2<sup>8</sup>



A hand-drawn diagram illustrating a grayscale image. It features a large square labeled "Gray Scale" at the bottom. Inside the square, there is a smaller square containing a stylized flower-like pattern. To the left of the main square, there is a vertical double-headed arrow labeled "h" above it and "!" below it, indicating height. To the right of the main square, there is a horizontal double-headed arrow labeled "w" above it and "!" below it, indicating width. Below the "Gray Scale" label, there is a red diagonal line with the text "24x24x1" written along it, representing the dimensions of the image as a 24x24 matrix with one channel.

24  
Correlation ↘

$$\frac{0 - 255}{\text{Intensity}}$$

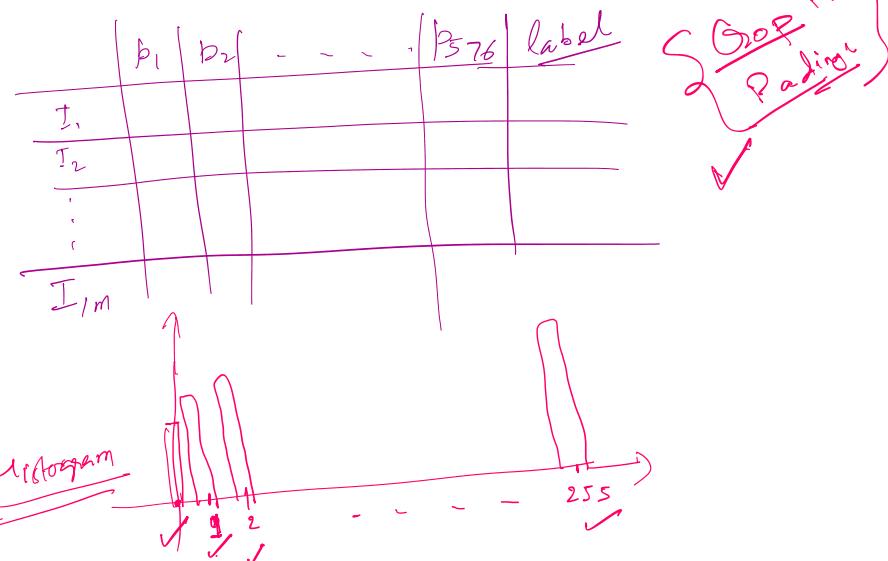
576

The diagram illustrates a mapping from a 4x4 matrix to a list of row vectors. On the left, a 4x4 matrix is shown with its first three rows labeled as row<sub>1</sub>, row<sub>2</sub>, and row<sub>3</sub>. An arrow points from this matrix to a list of row vectors on the right. The list contains four entries: [row<sub>1</sub>], [row<sub>2</sub>], [row<sub>3</sub>], and [row<sub>4</sub>].

$$I \in \mathbb{R}^{24 \times 24 \times 3} \rightarrow [0, 1, 200, 243, \dots] \in \mathbb{R}^{576}$$

$$I \in \mathbb{R}^{24 \times 24} \rightarrow [ ]$$

$$I^2 \in \mathbb{R}^{24 \times 24} \rightarrow [ ]$$



	0	1	2	...	255	label
$I_1$	.	.	.	.	.	
$I_2$	.	.	.	.	.	
:	.	.	.	.	.	
$I_m$	.	.	.	.	.	

ML model

Two types  $\begin{cases} \text{order of matter} \\ \text{order matter} \end{cases}$  Pixel values does not location of object do  
Pixel is important

$$\text{Language: } \sum = \{ w_1, w_{256} \} \quad \text{Dict: } S = \{ 256 \}$$

$$I_1: [0, 1, 200, 243, \dots] \quad \text{Text: } \underbrace{[ \dots ]}_{576}$$

$$I_2: [ ]$$

;

$$T : T$$

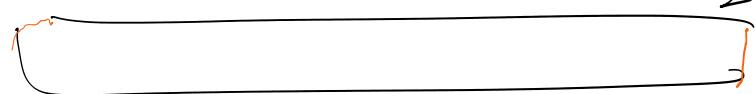
$\underbrace{[ \dots ]}_{200^n}$

~~one-hot-vector~~  
~~T~~: [  
~~W<sub>1</sub>~~ → [1000 0 - - - ]  
~~W<sub>2</sub>~~ → [0100 - - - ]  
~~256~~  
~~256~~  
 Assignment part

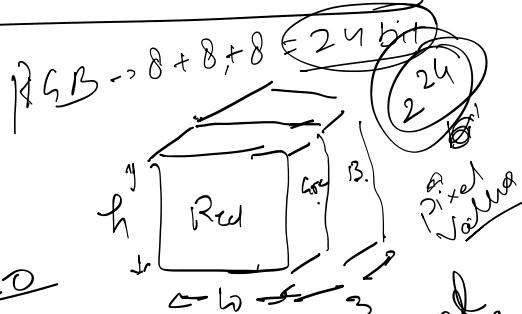
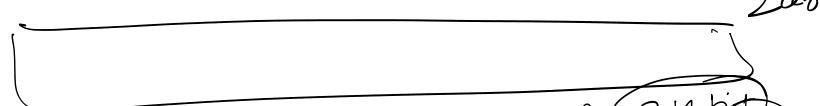
Text Data → Image Data  
Plot 2 View  
 Write some story about this

10x20

S1 -



S2



10x20

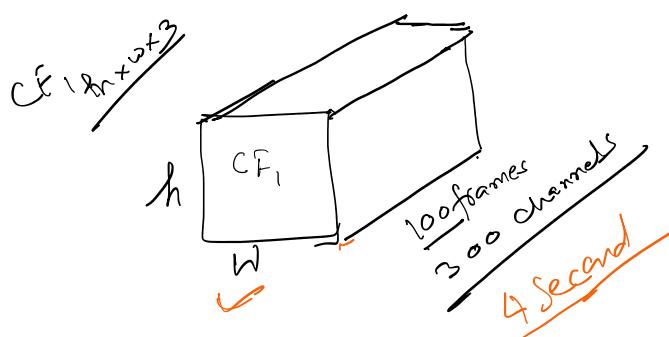
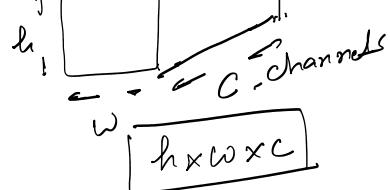
h x w x 3

channel

10 mins:

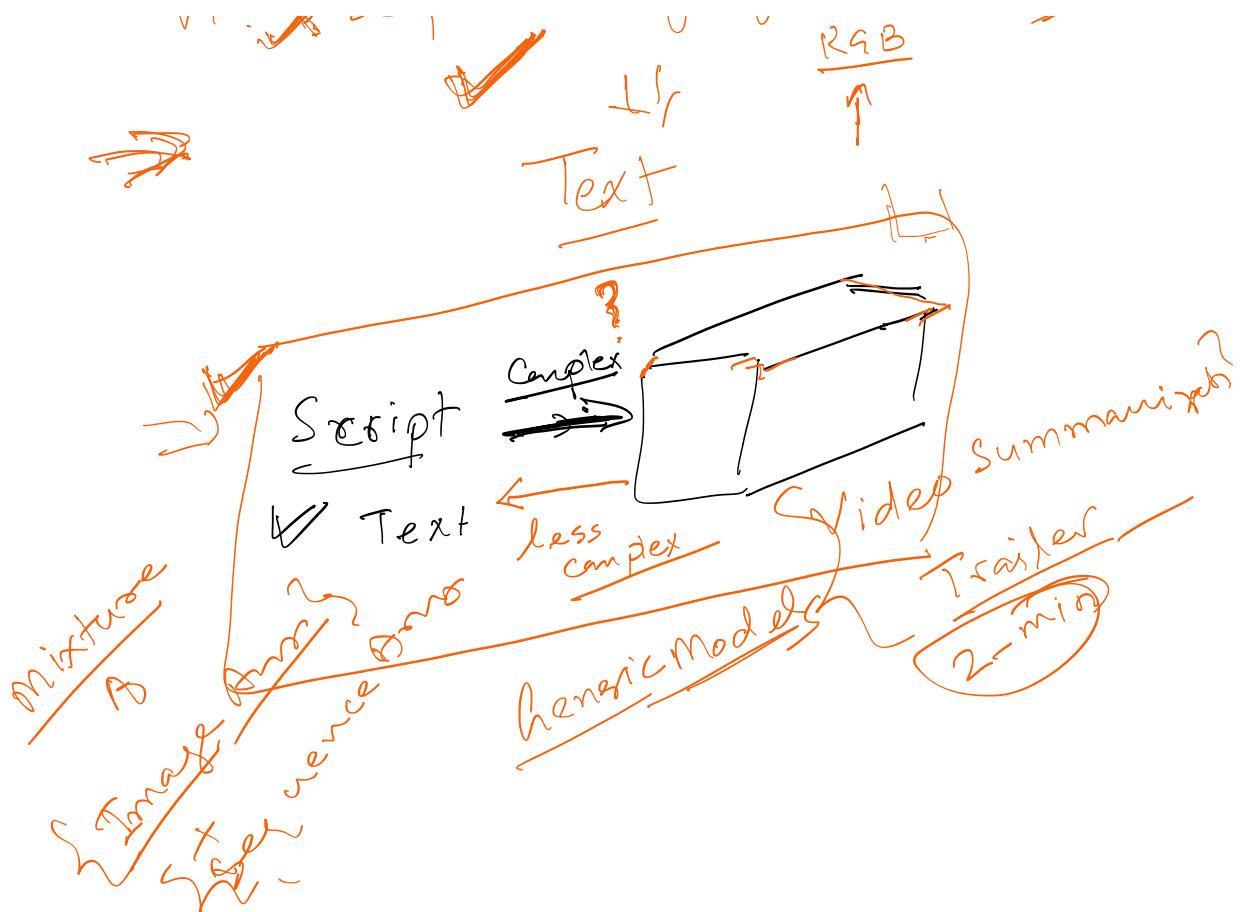
Video:

frame rate  
25 FPS



10 min video  
(10 x 60 x 25)  
frames

VI: ~~Sequence of frames~~   
 ✓ 1 hr RGB n ➤



## Text Data

Unstructured  $\rightarrow$  TF matrix

Unigram	$w_1$	$w_2$	$\dots$	$w_d$	lab. Senti.	T
$s_1$					+1	
$s_2$					+1	
$\vdots$					-1	
$s_m$					-1	

Sentences / Document

Information:  $w_i \rightarrow S_{w_i}$   $\text{fol.}$

The  $w_j \rightarrow S_{w_j}$   $\text{fol.}$

$m \rightarrow \text{unique.}$   $\sum_{i=1}^K s_{ij} = n$

$D = -\sum_{i=1}^K p_i \log p_i$

$D = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K p_{ij} \log p_{ij}$

The

$p_i \Rightarrow$  Unique word

$\sum_{i=1}^n$

$\log$

$w_i \rightarrow S_{w_i}$

# Document/Sentences =  $N$

IDF

$w_i \rightarrow S_{w_i}$

$S_{w_i} \over N$

$\log \left( \frac{N}{S_{w_i}} \right)$

{ How many documents/Sent  $(S_{w_i})$  contains  
the word  $w_i = S_{w_i}$  }

$S = \{S_1, \dots, S_N\}$

$k^{th}$  sentence:

$w_i \Rightarrow TF_{w_i} \times \log \left( \frac{N}{S_{w_i}} \right)$

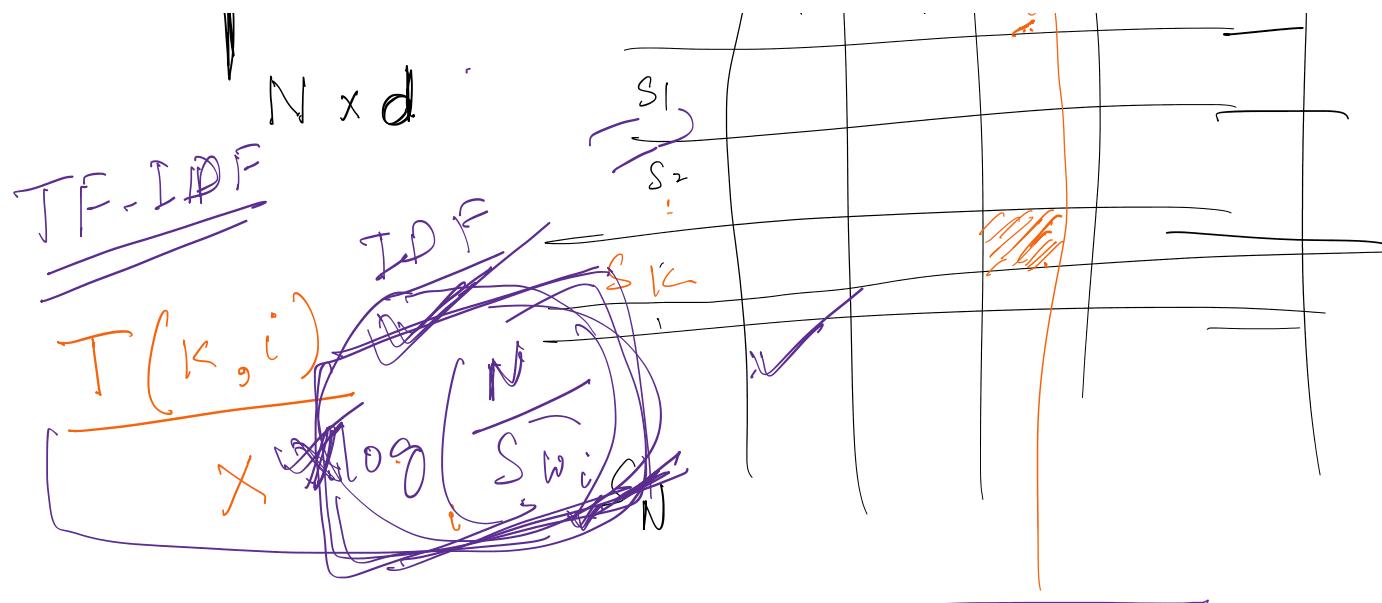
Term ~~for~~ Matrix

$T$

$N \times d$

$S_{w_i}$

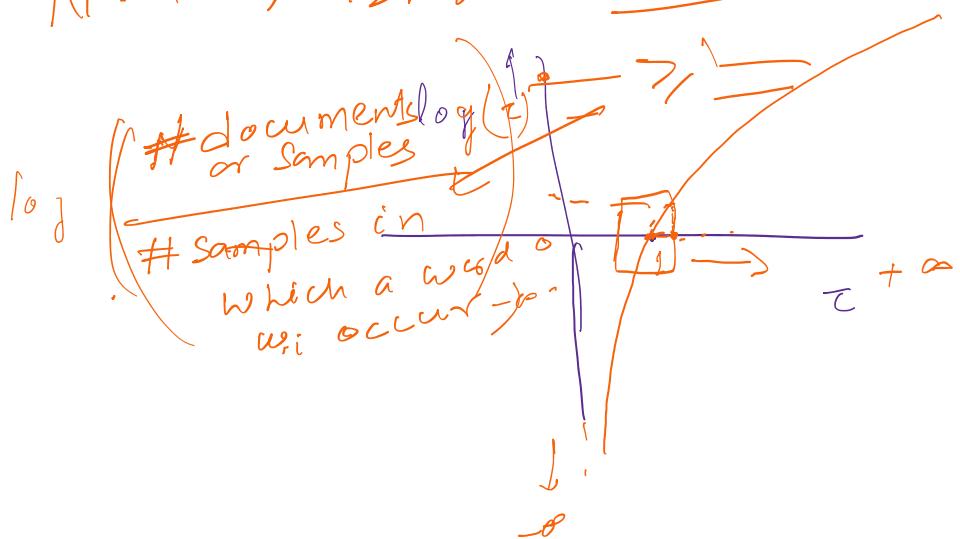
$S_{w_1}$	$S_{w_2}$	$S_{w_i}$	$\dots$	$S_{w_d}$
$w_1$	$w_2$	$w_i$	$\dots$	$w_d$
$S_1$				



• "The" 95%

$$\text{TF}(\text{the}) = 2 \times \log\left(\frac{100}{95}\right) \xrightarrow{1} \text{small +ve}$$

$$N = 1m, 95\% \text{ of } 1m = \underline{\underline{95 \times 10^6}}$$



### K-gram model

#### $K=2$ bi-gram model



• bi Windows Size = 1

S2: He is reading a book ✓

S3: Ramayana is a good book. ✓

(This, is) ✓  
(is, NLP) ✓  
(is, This) -

S2: He is reading a book ✓ (is, NLP) ✓

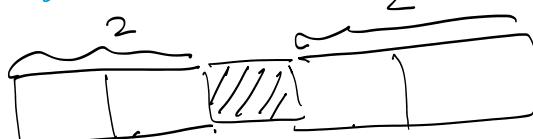
S3: Ramayana is a good book. ✓ (is, This) ✓ (NLP, class) ✓

S1: { (This, is), (is, NLP), (NLP, class) } (NLP, is) ✓  
1-ref (is, This), (NLP, is), (class, NLP) (class, word)

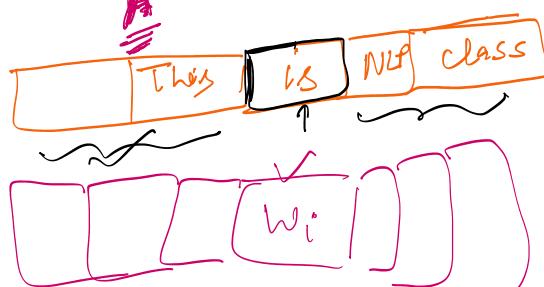
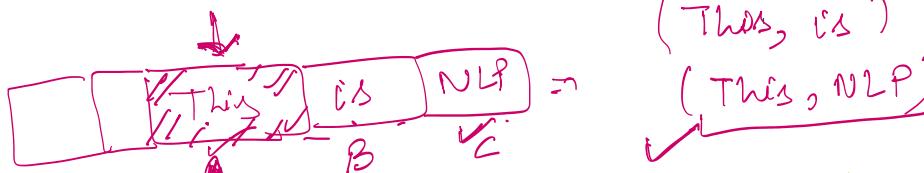
S2: { (He, is), (is, reading), (reading, a), (a, book)  
(is, He), (reading, is), (a, reading), (book, is) }

S3: { (R, is), (is, a), (a, good), (good, book)  
(is, R), (a, is), (good, a), (book, good) }

Window size = 2

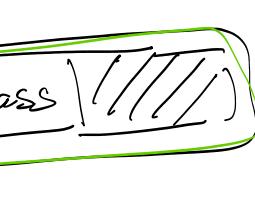


S1: This is NLP class. ✓



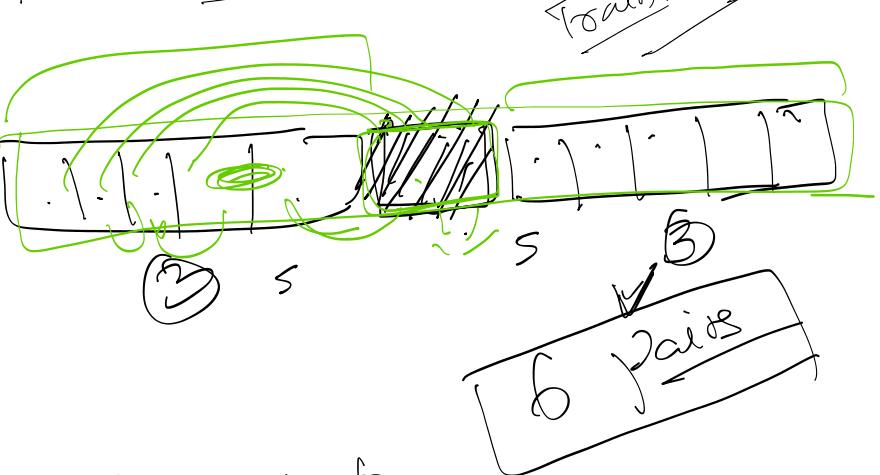
(This, is)  
(This, NLP)

(is, "NLP")  
(is, class)  
(is, This)  
(NLP, class)  
(NLP, is)  
(NLP, This)



Sample this Random Sampling

Training complexity



→ Combined few words to a single word  
↳ Phrases.  

Token = each phrase must be a single token

Unigrams x 2. Window Size

The diagram shows the mapping of words to their corresponding one-hot vectors. At the top, three words are listed: 'happy' (with a checkmark and a 100% label), 'heat' (with a checkmark and a 10! label), and 'm' (with a question mark). Below each word is its one-hot vector representation. The vector for 'happy' has a 1 at index 0 and 0s elsewhere. The vector for 'heat' has a 1 at index 1 and 0s elsewhere. The vector for 'm' has a 1 at index 2 and 0s elsewhere. Arrows indicate the mapping from the words to their vectors.

$[t]_{im} \quad [i]_m$

Word  
(Tokens)  $\rightarrow$  Vector

one-hot-vector

• shorter in length

+  
semantics  
or context

similarity  
b/w words  
tokens  
+  
semantically  
x context

$w_i \rightarrow$  5 closest words

$(w_i, \text{Dict}) \Rightarrow \text{prob}$

$(w_i, \text{Dict}) =$

$w_i \rightarrow 2 \checkmark \quad 2$   
 $(w_i, w_1) \rightarrow 3 \cdot 3 \checkmark w_i$   
 $(w_i, w_2) \rightarrow 0 \quad 0$   
 $(w_i, w_3) \rightarrow 0 \quad 0$

$(w_i, w_m) \rightarrow 1 \checkmark$

$(w_i, w_2)(w_i, w_3) \rightarrow \underline{.79}$

if  $(w_i, w_3) > .9$  &  $(w_i, w_2) > .9$

the  $w_i, w_j$  ~~are~~ they may  
semantically close

Labeled  
Assignment 1 Part D

1. TF-IDF Matrix

7

## Assignment 1 Part D

↳ TF-IDF Matrix

↳ Generate Bi-grams using skip-gram  
with window size 2-3.

↳ Generate & store in a file