

NLP

Next revolution will come when we are able to interact with machines the same way we interact with humans.

The functionality of search engines is powered by NLP. Other modern applications include

- 1) Contextual Advertisements
- 2) Social Media: removing adult content, opinion mining
- 3) Search Engines.
- 4) Chatbots. (let's say in Swiggy or Amazon : Customer Service)

Try thinking more use yourself!!

Approaches to NLP

1) Heuristic Approach "Jugaad" approach

Rule based approaches were created by NLP practitioners.

Eg: Let's say for sentiment analysis for a movie review, we count the number of positive words & no. of negative words.

- Example:
- 1) Regular Expressions
 - 2) Word Net

3) Open Mind Common Sense.

↳ Open ^{community} for everyone (open source)

↳ Anyone can contribute common English facts. Basically a database was created to be used for NLP tasks.

These are quick approaches, error free because humans are constructing it.

2) Machine Learning Approach

Good for open ended problems.

We cannot make general rules everywhere.

Here, our algorithms construct the internal pattern required from Input and Output.

ML Algos used:

- 1) Naive Bayes
- 2) SVM

3) Logistic Regression

4) LDA

5) Hidden Markov Models

Natural Language Processing is a branch of artificial Intelligence (AI) that helps machines to understand and process human languages either in text or audio form.

Challenges to NLP

- 1) Ambiguity in the language
- 2) Contextual words.
Eg: River bank & Money Bank
- 3) Colloquialisms & slang. Eg: Piece of cake
- 4) Synonyms.
- 5) Irony, Sarcasm & tonal diff'n
- 6) Spelling Errors.

Bottleneck of

→ Heuristic Approach.

3) Deep learning Approach

→ Automatic feature and parameter generation.

1) RNN

2) LSTM

3) GRU/CNN

4) Transformers

5) Auto encoders.

They don't care about sequential

data → which is fixed by RNN and

(DL) others.

NLP Pipeline

The following steps are necessary for constructing any NLP application

① Data Acquisition: Acquiring relevant data

Even if we use Heuristic Approach, we need extra data for better performance. (More data is always good)

2) Tent Preparation:

i) Text cleanup: fixing spelling mistakes, incomplete phrases

ii) Basic preprocessing: Tokenization of words.

iii) Advanced preprocessing
Eg: Parts of Speech, Chunking,
Tagging.

3) Feature Engineering: We need data in correct format

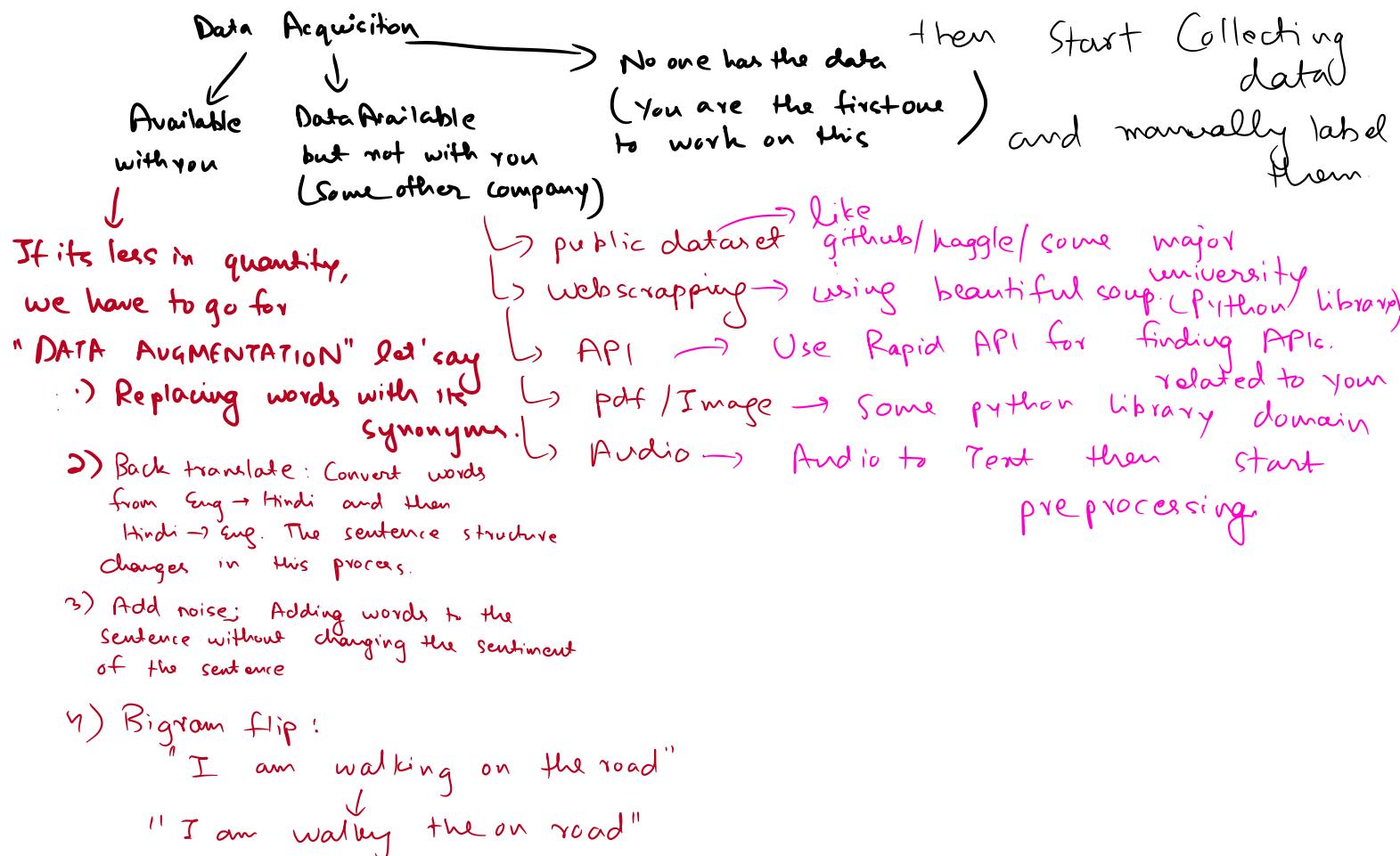
We need to change our text data into numerical

data · ~~Eg.~~ Word2Vec, TF-IDF, One-hot encoding.

4) Modeling: Actual Algorithm is applied

5) Deployment: Monitoring & Model Update

The above pipeline is not universal. For elementary NLP tasks like POS, Text Representation, this works fairly well.



Text Preparation

Cleaning

↳ HTML tag removal

↳ emoji from Social Media are removed

↳ Spelling checker

None of them are mandatory
↳ lib called Textblob is used
(task dependent)

Basic preprocessing:

Basic

→ Tokenization
(word & Sentence Tokenization)

Optimal

→ Stop word Removal.
(., ;, ,)

Don't have semantic meaning

→ Stemming (brings words to their root form)
(dance, dancing, danced)

All 3 of them have same semantic meaning

→ Removing digits, punctuations

→ Lowercasing

→ Language detection.

Advanced Preprocessing

→ POS tagging

→ Parsing

→ Coreference resolution

Feature Engineering

When we are using Machine learning approach

* We have to do the feature engg. ourselves

(Model's Interpretability is high)

Deep approach

* feature engg. is done by the neural networks.

(Model's Interpretability is low)

Modelling

Data is ready, now we have to fit a model to find patterns.

possible ways:

heuristic : if data is less, we try to fit a general rule to define trends. let's say we are making a spam email classifier

ML Algo : Data (more data)

DL Algo : Data (even more data) → use Transfer Learning.

(BERT)

We have a model trained on a very large dataset and it has large capabilities.

(let's say google auto complete)

evaluation

intrinsic
(model centric)

extrinsic
(business why product is centric not being used?? by)

Cloud API :

Decided by nature of data & Nature of Problem

for heuristic approach, we take a note of the email id from which we are continuously receiving spam emails

Text Preprocessing

Some of the common text preprocessing/cleaning steps:

- * Lower casing
 - * Removal of Punctuations
 - * Removal of Stop Words
 - * Removal of Frequent Words
 - * Removal of Rare Words
 - * Stemming
 - * Lemmatization
 - * Removal of emojis / emoticons
 - * Conversion of emoticons to words
- * Removal of URLs
 - * Removal of HTML tags
 - * Chat words conversion
 - * Spelling correction

Check Out

"Text-Preprocessing - AIWinter.ipynb"

Feature Engineering / Text Representation

We convert texts into numbers (numerical data)

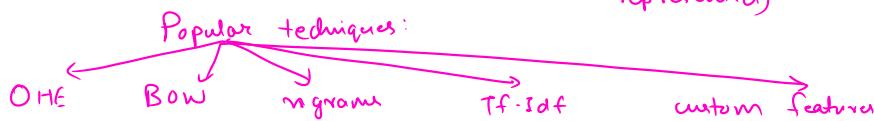


These are represented by vectors

So, this step is also called "text vectorization"

The numbers should capture the semantic meaning of the text

↓
(hidden meanings should be represented)



Terms:

1) Corpus (C): When we concatenate all of our strings
let's say we have movie reviews.

When we combine all movie reviews together (many words might be repeating)

2) Vocabulary (V): Unique words from corpus.

3) Document (D): each movie review

4) Word (w): Each word of document

word

Eg: Review 1: "I loved SRK in Jawaan" Sentiment 1

Review 2: "All actors had great chemistry" 1

Review 3: "Boring & predictable story line" 0

I haven't removed stop words here
Do remove if required

Corpus: I loved SRK in Jawaan All actors had great chemistry Boring & predictable story line. → 15 words.

Vocabulary (V): 15 words (unique words)

Documents (D): Rev 1, Rev 2, Rev 3

① One hot encoding Sentiment So, Corpus:

S1: IITK is my college	1	IITK is my college My College is IITK College is IITK	"stop words not removed you might choose to remove"
S2: My college is IITK	1	Vocabulary:	
S3: College is IITK	0	IITK, is, my, college	

we represent words in

$$\begin{array}{l} \text{HTRK} \quad \text{is} \quad \text{my} \quad \text{college} \\ \text{HTRK} = [1, 0, 0, 0] \\ \text{is} = [0, 1, 0, 0] \\ \text{my} = [0, 0, 1, 0] \\ \text{college} = [0, 0, 0, 1] \end{array}$$

We represent words in terms of vectors of length V .

$$\checkmark = 4$$

So, S1 is converted to $\left\{ [1, 0, 0, 0], [0, 1, 0, 0], [0, 0, 1, 0], [0, 0, 0, 1] \right\}$ 4×4
 S2 : $\left\{ [0, 0, 0, 1], [0, 1, 0, 0], [1, 0, 0, 0] \right\}$ 3×4

Very logical, intuitive, easy to code.

Flaws: 1) Sparsity

Let's say our vocabulary is 50000 (corpus of some million words)

So, HTP would be $[1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, \dots, 0, 0, 0, 0]$
 (Sparse array)

leads to overfitting (lots of zeroes in it)

2) Sentence is represented by Words & Vocabulary has to be kept same for any ML Model.
All sentences don't have same size. Input size

3) Out of Vocabulary (OOV) Problem:

Let's say we have to predict the sentiment of the sentence:

"IITS is better"
out of Vocabulary words So, we want predict the Sentiment.

Let's say:

walk, run, bottle = Vocabulary
verb = Representation

walk :	$[1, 0, 0]$
run :	$[0, 1, 0]$
bottle :	$[0, 0, 1]$

A diagram illustrating vector representation in 3D space. A point P is shown with three vectors originating from it: one labeled "run" pointing along the y -axis, one labeled "walk" pointing along the x -axis, and one labeled "bottle" pointing along the z -axis.

So, distance b/w sum-walk
= dist. b/w bottle walk.

But, run I walk have a
closer meaning compared
to the other pair.

So, the Semantic meaning is not being captured, which kills our purpose of vectorization.

2) Bag of Words (BoW)

Core Idea:

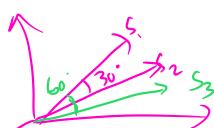
Sentences with same semantic meaning will have same type of words that appear in similar frequencies in the sentence.

Order of the words does not matter

content does not matter (but at the end its while writing the vectors captured indirectly)

All text is converted to N -dimensional vector.

Let's say $N=3$ sentence is given a vector in the N -dimensional



If the dot product is high, they are similar to each other.

I haven't removed stop words in the example
remove it if required

$$\text{Let's } \langle \text{Le bhn } S_1, S_2 : 30^\circ \rangle$$

$$\langle \text{Le bhn } S_1, S_3 : 60^\circ \rangle$$

$$S_1 \cdot S_2 = \cos 30^\circ = \frac{\sqrt{3}}{2} \quad S_1, S_2: \text{more similar}$$

$$S_1 \cdot S_3 = \cos 60^\circ = \frac{1}{2} \quad (\text{Because they are close to each other.})$$

$\langle \text{Le } \rangle$ is lesser.

Scikit Learn has

{ count - vec }
functionality

max-features

binary
{ binary = false;
if true
then freq > 1 will be made }

Sentence [1, 0, 1, 2, 3] \rightarrow [1, 0, 1, 1, 1]
vector

	1	2	3	4	5
Boy	1	1	1	1	0
Bauder is a boy	0	2	1	0	0
Bauder is Bauder	1	1	0	0	1

S_1 : Bauder is a boy

S_2 : Bauder is Bauder

S_3 : Boys hate Bauder

Vocabulary = 5

So, 5^(v) dimensional vectors for each sentence

max-features: when many rare words are present. It decides

the number of words considered for analysis
max-features = 1 So, top freq. words are

max-freq = 2 So, top 2 freq. words observed are analyzed.

Bottlenecks:

1) Sparsity: The dim. of vectors will be very large (let's say 50,000)

So, many zeroes will be present, end up overfitting

2) OOV: Improvement from OIE
This technique won't report an error.

3) Ordering: We are dealing with frequencies

4) S_1 : Bauder is a very good singer and a nice boy
 S_2 : Bauder is not a very good singer and a nice girl

Vectors might be close but semantic meaning is very different.



3 Bag of n-grams

Single word were creating issues

'not good' and 'good' have different meanings

$n=2$ (Bigram)

We take pairs.

Example

S_1 : Baruden is your mentor

S_2 : Mani is your mentor

S_3 : khurhal is your mentor and mani is your mentor.

I haven't removed stop words here
, remove it if required

Tokens

Vocab:

Baruden is | is your | your mentor (Mani is | khurhal is | mentor and) and Mani

S_1 :				0	0	0	0
S_2 :	0			1	0	0	0
S_3 :	0			0)	1	1

We used bigram that's why we can create greater angle b/w 'not good' & 'good'
If we use trigrams, difference would be more clear.

So, Semantic meaning is captured better of the sentence

Benefits

- 1) Easy to implement
- 2) Able to capture

Disadvantages

- 1) We still ignore OOV
- 2) Uni-gram, bi-gram, tri-gram } dimensionality ↑ as tokens increase
So, this increases the model training time.

Term

Inverse-document frequency

4 Tf - Idf

Term Frequency

Inverse-document frequency

Core Intuition:

If a word appears (rare) less in the corpus, then its weight increases in the sentence it's present in.

S_1 : Good old Baruden boy

$$TF(\text{Baruden}) = \frac{1}{6}, IDF(S_1) = \ln\left(\frac{3}{3}\right) = 0$$

S_2 : Baruden nice boy

S_3 : Old Baruden good boy

$$TF(\text{nice}) = 0, IDF(\text{nice}) = \ln\left(\frac{1}{1}\right) = \ln 1$$

nicer is rare in the corpus.

(Rare words will have higher IDF)

$$TF(t, d) = \frac{\text{No. of occurrences of term t in doc d}}{\text{Total no. of terms in doc d}}$$

$0 < TF < 1$ (Probability) { If $TF \uparrow$: So, term t is more frequent }

$$IDF = \ln \left(\frac{\text{Total no. of docs in corpus}}{\text{No. of docs with term t in them}} \right)$$

No. of docs with term t in them

Why \ln is used?

$$Idf = \left(\frac{N}{n}\right) ??.$$

$$\text{If } N=10000 \\ n=2$$

then $\frac{N}{n}=5000$ value is very large

So, when we are considering $IDF * TF$,
larger IDF values will dominate over TF values
We want both of them to contribute equally.

$$\text{So, } IDF = \log_e\left(\frac{N}{n}\right) = \ln(5000)$$

↳ Here, \log_e normalizes the value

$$\text{Sometime, } IDF = \ln\left(\frac{N}{n}\right) + 1 \quad [\text{In sklearn library}]$$

If any word is appearing in every sentence; its

$$IDF = \ln(1) + 1 = 1 \quad \text{instead of } (0)$$

Advantages

1) Used for Information Retrieval

Search Engines go through millions of webpages and give suggestions to you.

Disadvantages

1) Sparsity (High dimensions, so lots of zeroes)

2) OOV

3) Semantic Meaning is still not captured in some cases.

"Einstein was genius"

"Einstein was brilliant"

They kinda have similar meanings. But TF-IDF treats them differently.

Miscellaneous Custom features

$$\frac{\text{No. of Positive Reviews}}{\text{No. of Negative Reviews}}$$

Here, domain knowledge plays an important role.

TF-IDF Calc for SI

Each term of word Vocab for corpus!
in vector is represented by Good, old, Barudew, boy, nice
by TF * IDF

$$TF(\text{Good}) = \frac{1}{4} \cdot IDF(\text{Good}) \\ = \ln\left(\frac{3}{2}\right)$$

$$TF(\text{old}) = \frac{1}{4}, IDF = \ln\left(\frac{3}{2}\right)$$

$$TF(\text{Barudew}) = \frac{1}{4}, IDF = \ln\left(\frac{3}{3}\right)$$

$$TF(\text{boy}) = \frac{1}{4}, IDF = \ln\left(\frac{3}{3}\right)$$

$$TF(\text{nice}) = 0, IDF = \ln(3)$$

So, SI is represented by

$$= \left[\frac{1}{4} \ln\frac{3}{2}, \frac{1}{4} \ln\frac{3}{2}, \frac{1}{4} \ln 1, \frac{1}{4} \ln 1 \right]$$

$0 \times \ln 3$

$$= \left[\frac{\ln 1.5}{4}, \frac{\ln 1.5}{4}, 0, 0, 0 \right]$$

$$= [0.101, 0.101, 0, 0, 0]$$

Word Embeddings

- frequency types
 - BOW
 - Tf-IDf
 - Glove

& prediction types

- Word2Vec

Word2Vec

- 1) Semantic Meaning is captured
 - 2) low dimensional vectors are formed (300-500)
 - 3) dense vectors (non-zero elements)
 - ↳ overfitting
- Representing the Vector in 300 dimensional space.
- Vector Arithmetic works here. (Addition, Subtraction, Dot Product)
 - ↳ Dot product also works

We can use pre-trained as well as self-trained word2vec models
 ↓
 training it on our own dataset

	king	queen	Man	Woman	Monkey	
gender f_1	1	0	1	0	1	We can never find out what they mean exactly
Wealth f_2	1	1	0.3	0.3	0	
Power f_3	1	0.7	0.2	0.2	0	
weight f_4	0.8	0.4	0.6	0.5	0.3	
speak f_5	1	1	1	1	0	

$$\text{king} = [1, 1, 1, 0.8, 1]$$

Internally, Word2Vec creates features based on the vocabulary.
 These features might be different for different vocabulary.

$$\text{king} - \text{man} + \text{woman} \approx \text{queen}$$

We use automated techniques to create features.

We might have a vocabulary of 50 lakh words. So, manually generating features is quite tough.

So, we use Neural Networks.

All 300 numbers represent unique features, but we can never find out what they represent exactly.

Core Intuition

The underlying assumption of Word2Vec is that two words sharing similar contexts also share similar meaning and consequently a similar vector representation from the model.

Types of Word2Vec:

CBoW → skip-gram

Continuous Bag of Words (CBoW)

let's say: "Choose Architecting Intelligence for deep learning"
 we have.

We want to write the word embeddings for each word.
 (vectors)

Let's first decide the size of these vectors for each word (dimension)

(let it be 3)

$$\text{So, choose} = [-, -, -]$$

We create a dummy problem and try to solve it.
 Now, window keeps on moving for deep learning → Sentence
 Choose Architected Intelligence for deep learning
 We start by considering windows of size=3 [It can be 5, 7, 9, ... also]

1st window
 central word: target word. (Architecting)
 left & right words: context words (choose, Intelligence)

So, X (Input) | Y (target)
 Choose, Intelligence | Architecting

We created a "fill in the blanks" problem "Choose _____ Intelligence".
 It is Supervised ML
 By taking content words as inputs, we are trying to predict the target word.

Now we take OHE for our vocabulary

Choose: [1, 0, 0, 0, 0, 0]

Architect: [0, 1, 0, 0, 0, 0]

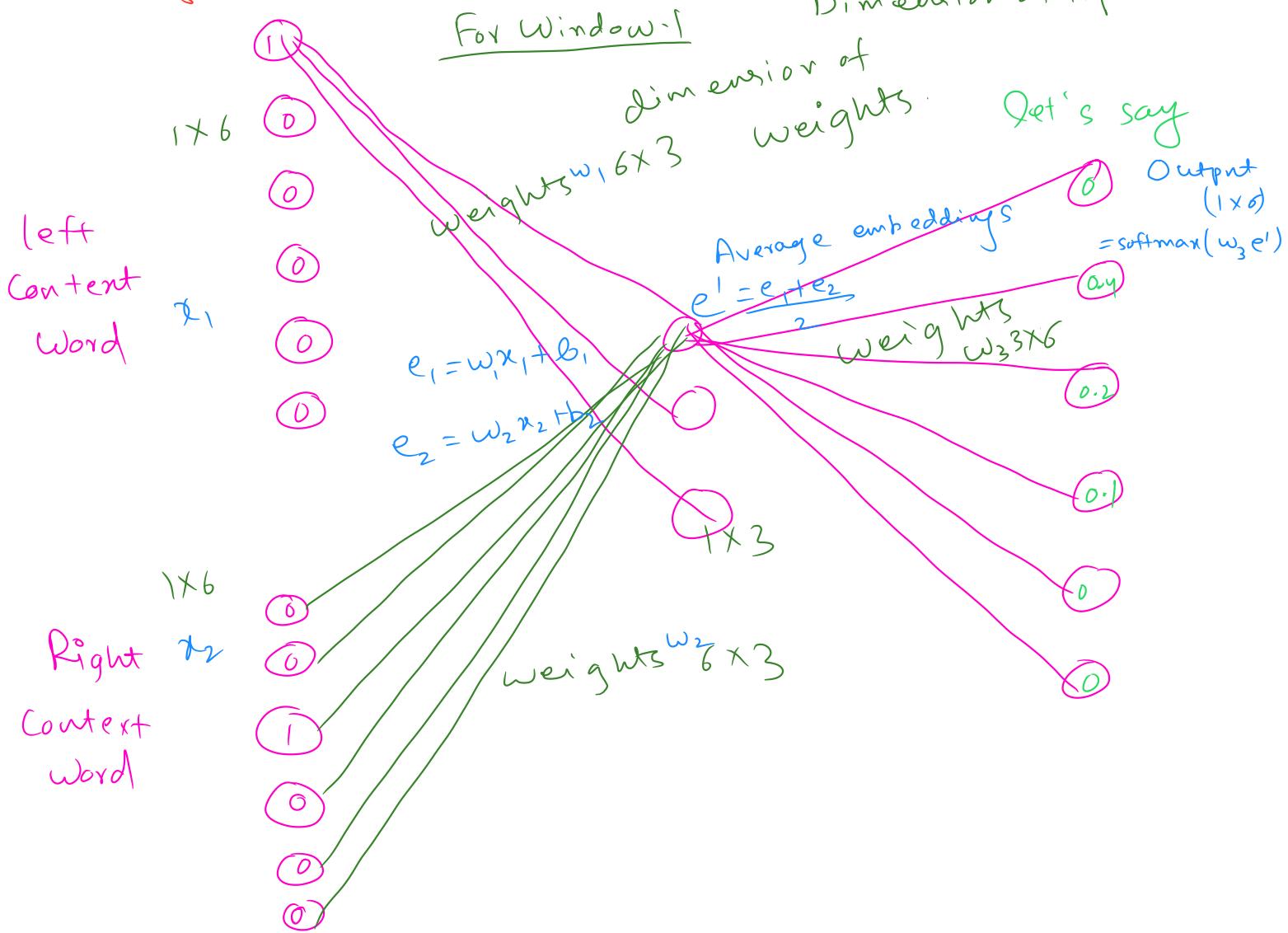
Intelligence: [0, 0, 1, 0, 0, 0]

for: [0, 0, 0, 1, 0, 0]

Deep: [0, 0, 0, 0, 1, 0]

Learning: [0, 0, 0, 0, 0, 1]

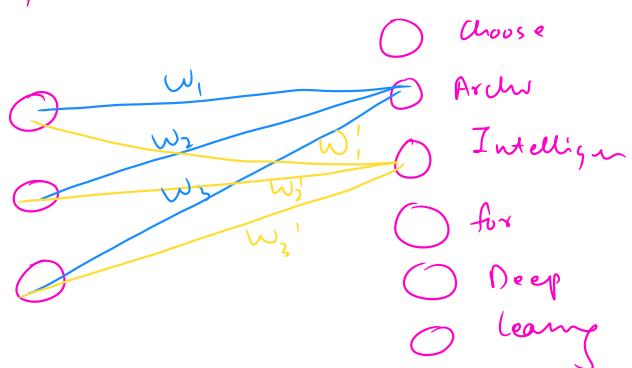
Dimension of input: 1x6



Expected output $\rightarrow [0, 1, 0, 0, 0, 0]$
 Output we got $\rightarrow [0, 0.4, 0.2, 0.1, 0, 0]$

So, we calculate loss and do back propagation

So, vector representation for words:



After back propagation is complete

Architecture: $[w_1, w_2, w_3]$

We wanted 3 dimensions
representation initially

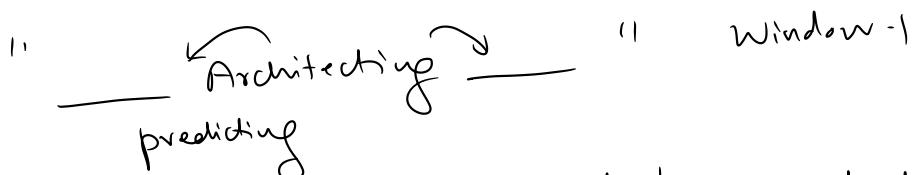
Intelligence: $[w_1', w_2', w_3']$

Skip gram

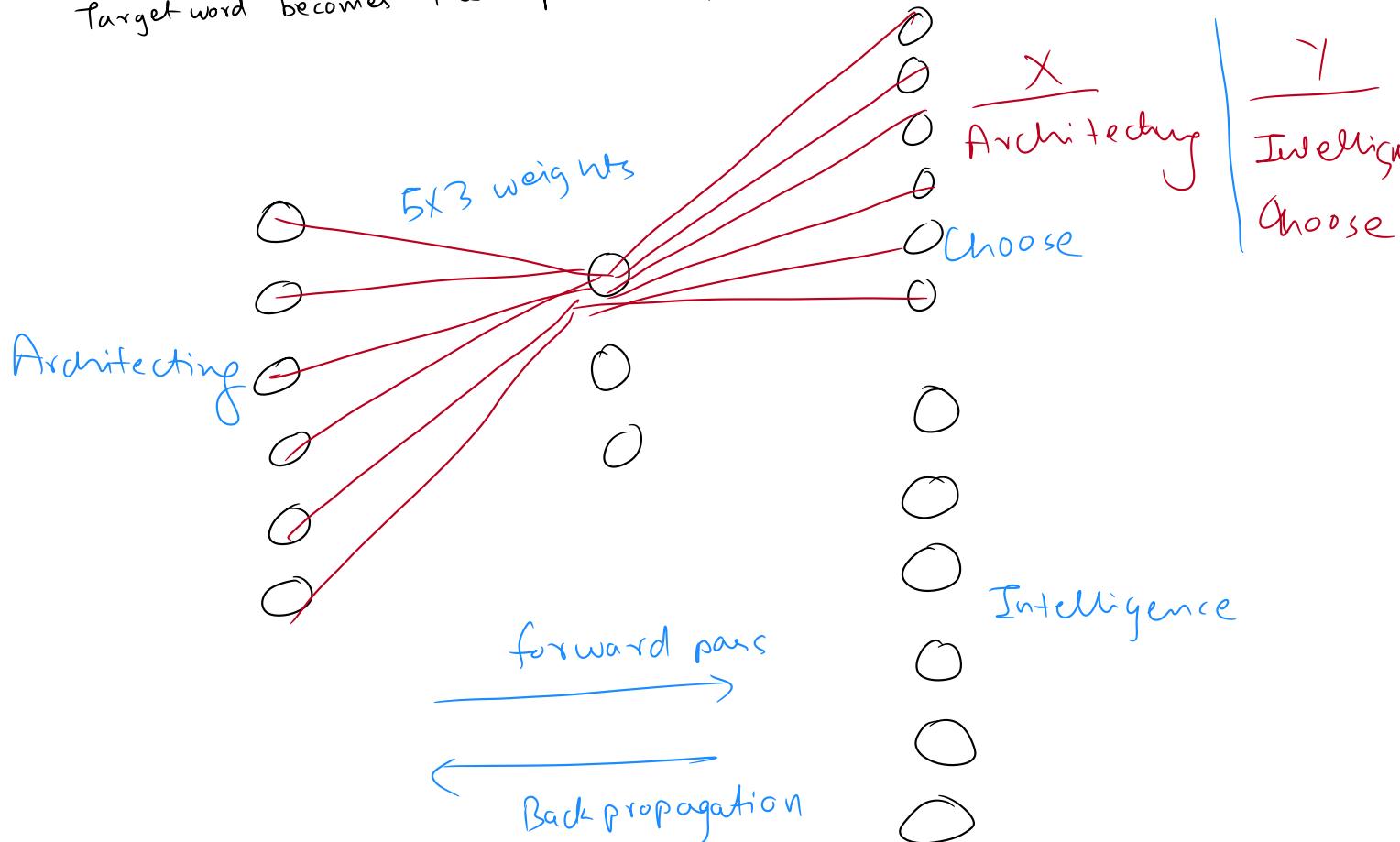
"Choose Architecting Intelligence" for deep learning"

Window-size = 3

Reverse our problem this time



Target word becomes the input here, content words become the output



For Small data: use CBOW
Large data : Skip gram

How can we improve Word2Vec?

- Increase the training
- Increase the dimensions of final vector embeddings (The Hidden layer)
- Increase window size.

New topic

Text Classification

Classification: Supervised Machine learning task

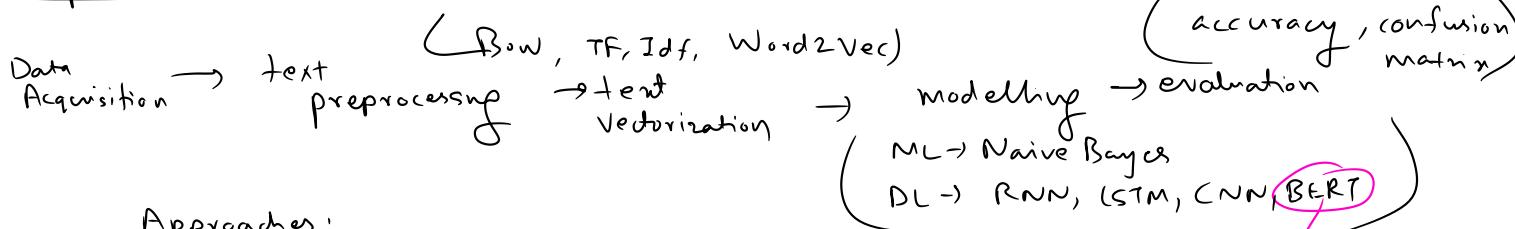
- ↳ On textual data: Email Spam Classifier
Movie Review Sentiment Analysis.

Types of Classification:

Binary	Multi Class	Multilabel
Eg Entertainment, business.		

Single news can be part of multiple sources.
 let's say
 SRK
 Bollywood
 Entertainment.

Pipeline:



Approaches:

Heuristic

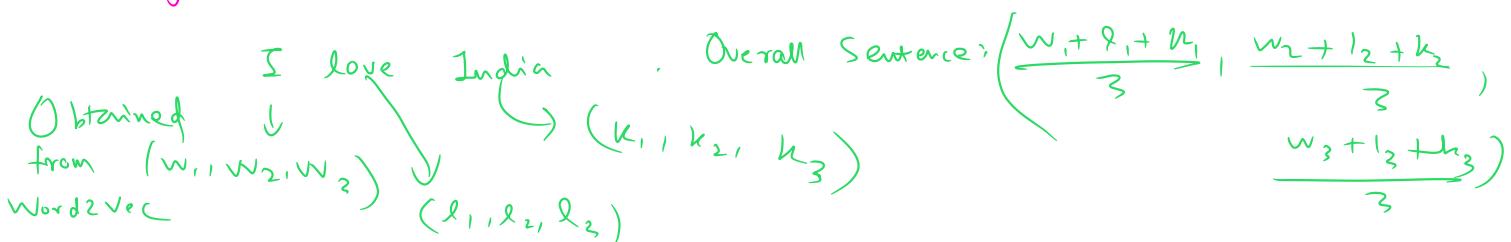
APIs (Google, AWS)

ML: Bow, TfIdf, n-gram

DL: word2vec

Word2Vec gives vector embeddings for each word

Avg Word2Vec: for sentence



General Guidelines:

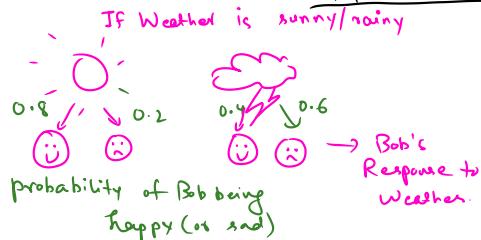
- 1) Ensemble techniques instead of Naive Bayes
- 2) Heuristic features should be combined with ML
- 3) Start ML techniques, then go for DL
- 4) Avoid Imbalanced dataset.

(POS) Parts of Speech Tagging

In simple words, we can say POS tagging is a task of labelling each word in a sentence with its appropriate part of speech.

POS tagging is a preprocessing step.
Used in Named Entity Recognition.
Question Answering System.
Word sense disambiguation
Chatbots.

Check out my Jupyter notebook! HIDDEN MARKOV MODEL



If Bob says that he was

Mon Tue Wed Thurs Fri Sat

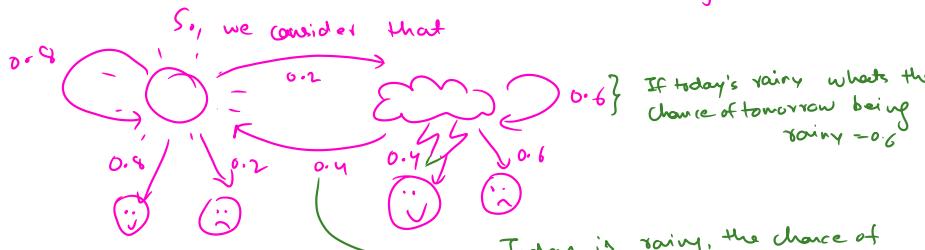


Based on Bob's review we have to guess the weather. So, it can be

S R S R S R

S: Sunny
R: Rainy

But the weather can't change so much.

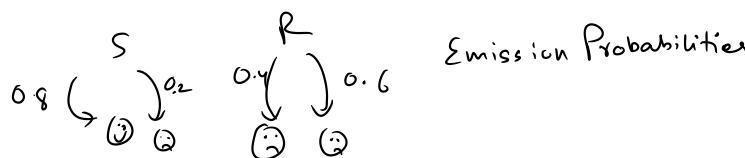
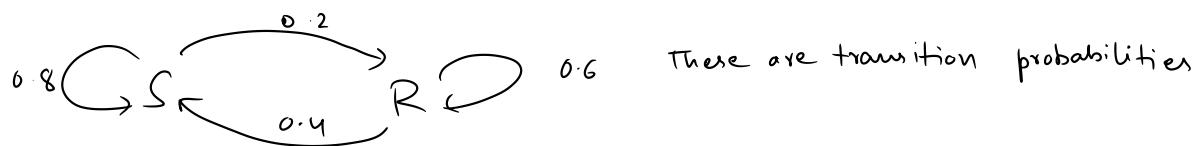


We can only see Bob's reaction, we can't see the weather.

[Bob talks to us over phone and tells us his mood. Bob leaves in France, we live in India].

So, Bob's Reaction: Observations.
Weather: Hidden States.

And its called the "Hidden Markov Model"



Q1) How we found these probabilities.
let's say the data of weather:

S → S → S → R → R → R → S → S → S → S → R → R → S → S → S

How many times S goes to R on the next day?

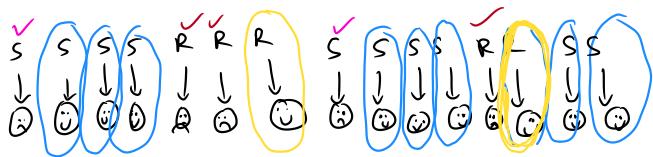
$$\begin{aligned} S \rightarrow R &= 2 \\ \text{Similarly, } S \rightarrow S &= 8 \\ R \rightarrow R &= 3 \\ R \rightarrow S &= 2 \end{aligned}$$

So, from S

$$\begin{aligned} S \rightarrow R &= 2 \\ 2+8 &= 10 \\ \frac{2}{10} &= 0.2 \end{aligned}$$

$$\begin{aligned} R \rightarrow R &= 3 \\ 2+3 &= 5 \\ \frac{3}{5} &= 0.6 \end{aligned}$$

Weather to Mood Data:



$$\text{So, } P(\text{☺}) = \frac{8}{8+2} = 0.8$$

$$P(\text{☹}) = \frac{2}{8+2} = 0.2$$

$$P(\text{☺}) = \frac{8}{15} = \frac{2}{3}$$

$$P(\text{☹}) = \frac{2}{15} = \frac{1}{3}$$

$$P(\text{S}) = \frac{10}{15} = \frac{2}{3}$$

$$P(\text{R}) = \frac{5}{15} = \frac{1}{3}$$

Q2) What is the probability that Random day is Sunny Rainy

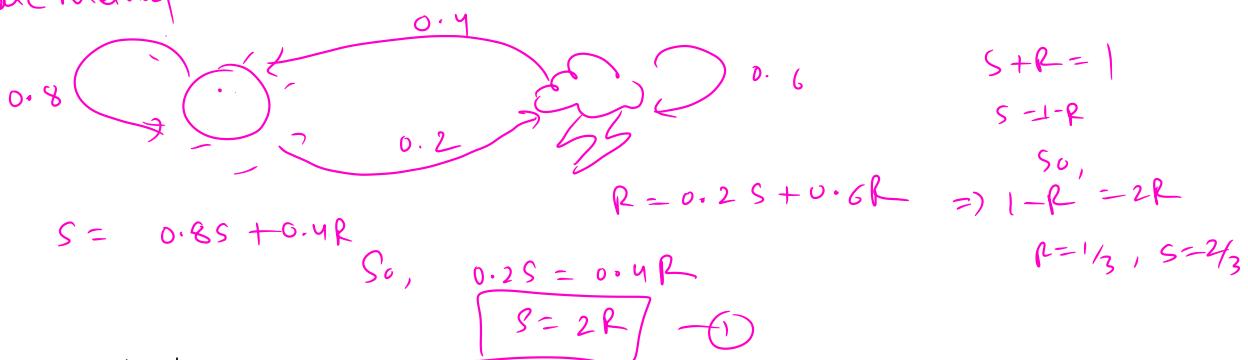
No. of S days = 10

No. of R days = 5

(If Bob gave us no inputs).

$$P(S) = \frac{10}{15} = \frac{2}{3}, \quad P(R) = \frac{5}{15} = \frac{1}{3}$$

Alternate Method



Q3) If Bob says "Happy ☺"

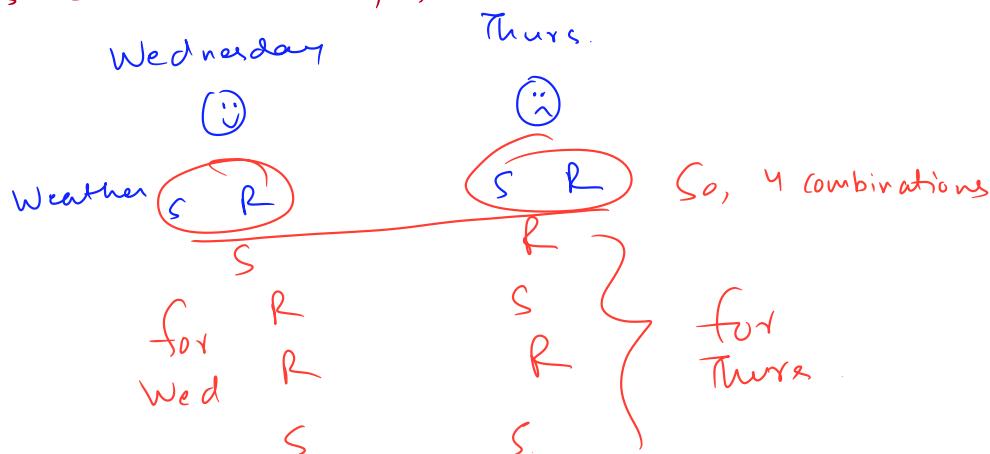
So, Probability of Sunny ↑↑ from 0.67

$$\text{So, } P(\text{Sunny}) = P(\text{☺}) \cdot \frac{P(\text{sunny})}{P(\text{☺})}$$

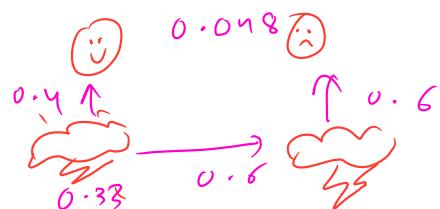
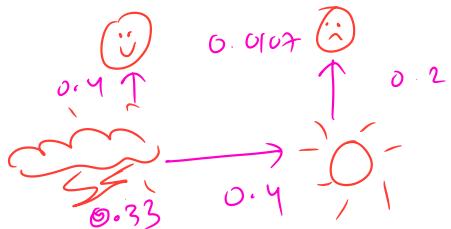
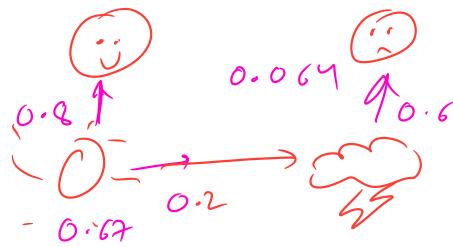
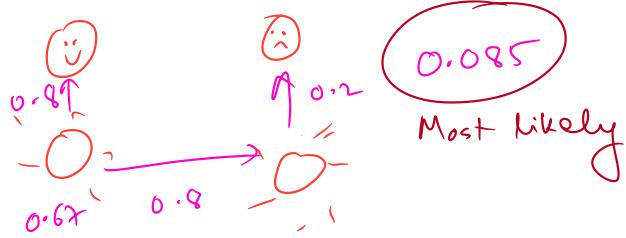
$$= (0.8) \times \left(\frac{10/15}{10/15} \right) = 0.8$$

(Bayes Theorem)

Q4) If for three days, Bob is ☺ ☹ ☺, what was the weather??
So, we have to guess the weather on all three days.
Let's start with 2 days;



Scenarios



So, Wed : Sunny

Thurs: Sunny

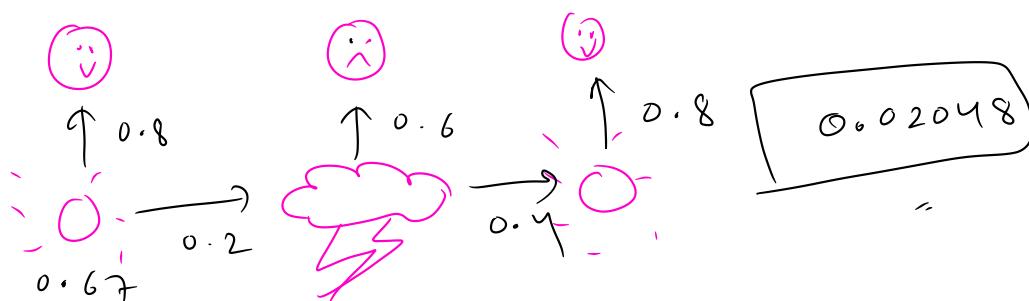
For three days

Wed	Thurs	Fri
8 combinations		
S	S	S
S	R	S
S	S	R
R	R	R
R	S	S
R	S	R
R	R	S
S	R	R

for



Combination



Check all 8 combinations



Best combination
(SSS)

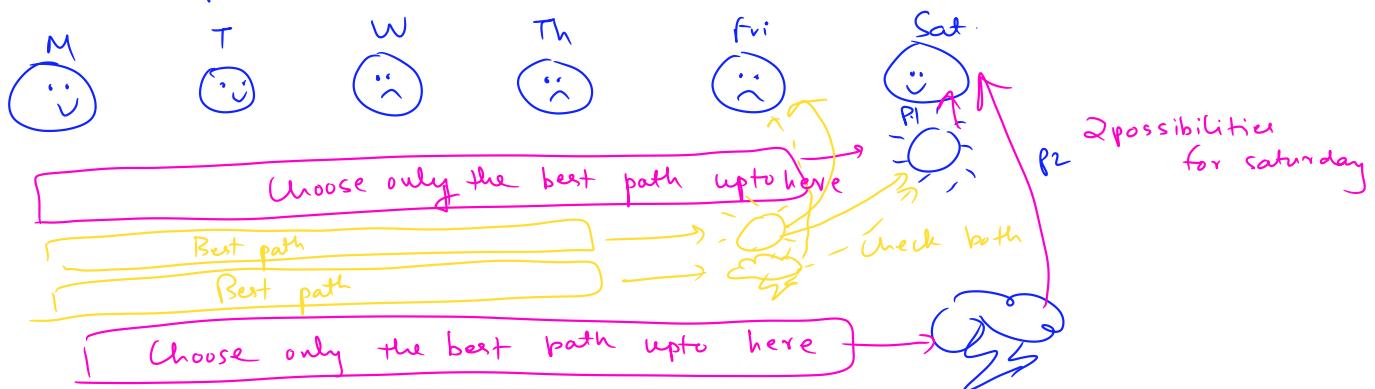
for 4 days, 16 (2^4 combinations)

The number of combinations grows exponentially (2^n).

So, to optimize, we have

VITERBI ALGORITHM

We need to optimize our solution



Check P_1 & P_2 and find out the better probability

POS tagging

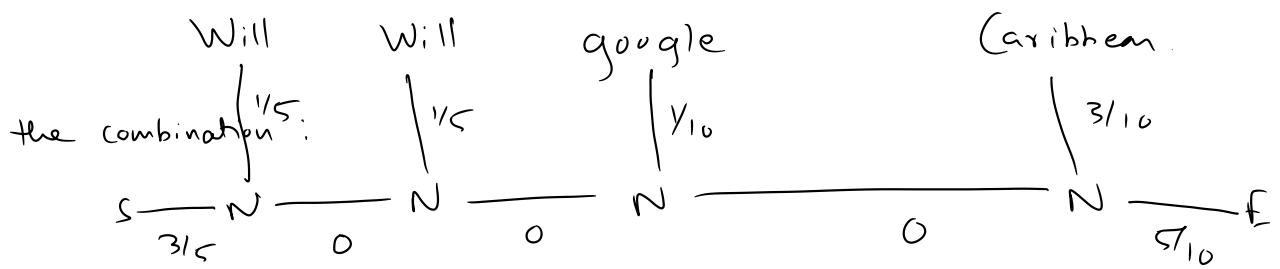
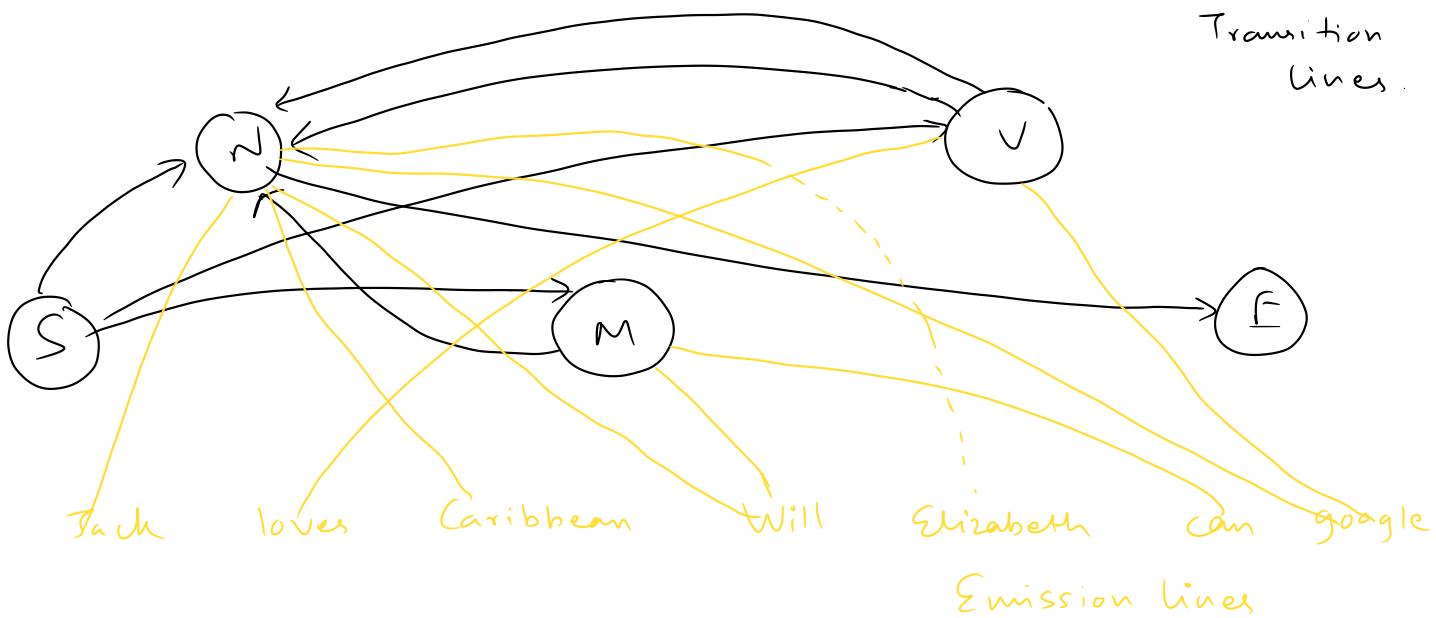
POS on the Sentence

Data		" will will google Caribbean"	
S	Jack	N	✓
S	lover	✓	Caribbean
S	Can	M	Jack
S	Jack	N	google
S	Will	M	Elizabeth
S	Elizabeth	N	google
S	Will	N	Caribbean
S	Elizabeth	N	loves
S	Will	N	Will
S	lover	✓	google.
S	Will	✓	E
S	lover	✓	N
S	Will	✓	E
S	lover	✓	M
S	Will	✓	V
			P(Jack Noun) = $\frac{2}{2+3+1+2+2} = \frac{2}{10} = 0.2$
Jack	(2) $\frac{1}{10}$	$0/2$	0
lover	$0/10$	$0/2$	$3/5$
Caribbean	$3/10$	$0/2$	0
google	$1/10$	$0/2$	$2/5$
will	$2/10$	$1/2$	0
Elizabeth	$2/10$	$0/2$	0
Can	$0/10$	$1/2$	0

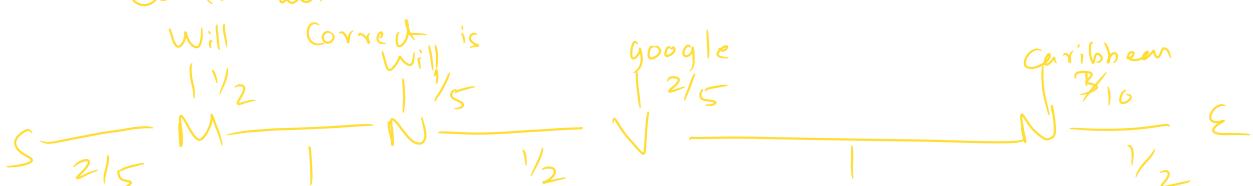
(Emission Probabilities are obtained)

For Transition Probabilities

	N	M	V	F
S	3/5	2/5	0	0
N	0/10	0/10	5	5
M	2/2	0	0	0
V	5/5	0	0	0



Check all :



Probability of best combination

$$\begin{aligned}
 &= 0.4 * 0.5 * 0.2 * 0.4 * 0.5 * 0.3 * 0.5 \\
 &= 0.012
 \end{aligned}$$

So, total possible combinations = $(3^6) = (81)$

But no. of POS $\gg 3$

So, its optimized by Viterbi

