

HMR)

- Eg. <s> martin justin can watch will <e>
 <s> spot will watch martin <e>
 <s> will justin spot martin <e>
 <s> martin will pat spot <e>

n: noun, M: model verb, V: verb.

→ ① Create Emission table . word : tag

<s>	martin	justin	can	watch	will	<e>
	N.	N.	M	V.	N.	.
<s>	spot	will	watch	martin		<e>
	N	M	V.	N.		y
<s>	will	justin	spot	martin		<e>
	M	N.	V.	N.		-
<s>	martin	will	pat	Spot		<e>
	N.	M.	V.	N.		.

	N	M	V	
martin	4			
justin	(2)			
will	1	3		
Spot	2		1	
can				
watch				
pat			2	
Σ	(9)	4	4	

$$P(\text{justin} | n) = \frac{2}{9}.$$

② Create state Trans.

Tag : Tag .

	<s>	N	M	V	<e>	
<s>	X (3)	1			X	$\therefore \# <s> = 4$
N	X	1	3	1	4	$\therefore \# N = 9$
M	X	1		3		$\therefore \# M = 4$
V	X	4				$\therefore \# V = 4$
<e>	X	X	X	X	X	
$P(n s) = s, n = \frac{3}{4}$						

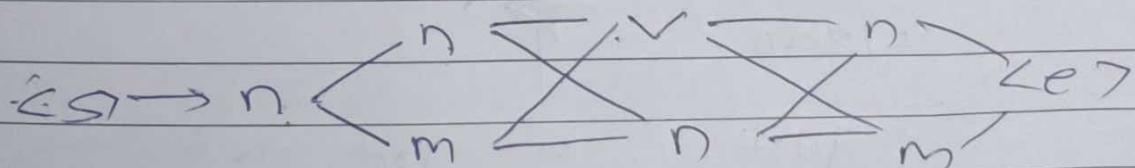
③ Test Data

<s> justin will spot will <e>

justin will spot will <e>
 n n v n n m ↗ Ambiguity

1 X 2 X 2 X 2 n

∴ Total ways : 8 .



$$\begin{aligned}
 & \begin{matrix} s, n \\ 1, 2 \\ 3 \end{matrix} \xrightarrow{\text{justin}} \rightarrow P(\text{justin}|n) * P(n|s) \\
 & \qquad\qquad\qquad = \frac{3}{9} * \frac{3}{4} = \frac{1}{6} .
 \end{aligned}$$

will as MV(m)

will as noun

$$\begin{aligned}
 & \begin{matrix} n, m, \text{will} \\ 1, 2, 3 \end{matrix} = (\text{will}|m) + (n|m) \\
 & \qquad\qquad\qquad = \frac{3}{9} + \frac{3}{9} = \frac{1}{4} \\
 & \qquad\qquad\qquad = \frac{1}{4} * \frac{1}{6} = \frac{1}{24} \checkmark
 \end{aligned}
 \qquad\qquad\qquad
 \begin{aligned}
 & \begin{matrix} a, n, \text{will} \\ a, n, n \end{matrix} = (\text{will}|n) + (n,n) \\
 & \qquad\qquad\qquad = \frac{1}{9} * \frac{1}{9} = \frac{1}{81} \\
 & \qquad\qquad\qquad = \frac{1}{6} * \frac{1}{81} = \frac{1}{486}
 \end{aligned}$$

∴ will becomes modal verb

Spot is a verb

∴ m → v → spot

$$\frac{3}{16}$$

↓

$$\frac{1}{128}$$

Spot is a noun

m → N → spot = Spot[n + fn]

$$= \frac{2}{9} + \frac{1}{9}$$

$$\frac{1}{18}$$

↓

$$\frac{1}{432}$$

will is a MV

v → m → will

$$\frac{3}{4} * (v, m)$$

(aplace smoothing)

$$\therefore \frac{3}{4} + \frac{1}{1000} = \frac{3}{4080}$$

$$\frac{1}{512000}$$

will is a n

v → n → will

$$\frac{1}{9} * \frac{1}{2} = \frac{1}{18}$$

$$\frac{1}{1152}$$

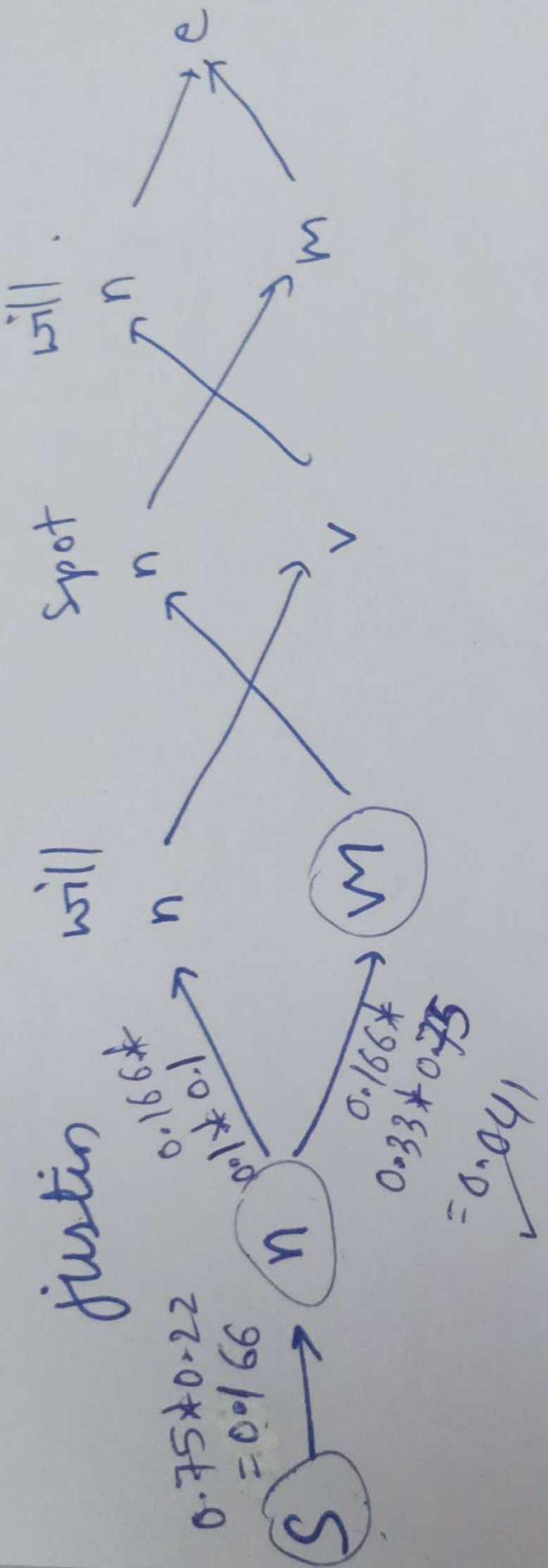
justin will spot will

n

m

v

n



Max. Entropy model

exp means e raised to

$$p(\text{spam} | x) = \frac{\text{spam}}{\text{total}} = \frac{\exp(0.6 + 0.7 + 0.5)}{\exp(0.6 + 0.7 + 0.5) + \exp(-0.3 - 0.4 - 0.2)}$$
$$= \frac{\exp(0.75)}{\exp(0.75) + \exp(-0.9)}$$

higher

$$p(\text{non spam} | x) = \frac{\text{non spam}}{\text{total}} = \frac{\exp(-0.3 - 0.4 - 0.2)}{\exp(0.75) + \exp(-0.9)}$$
$$= 0.25.$$

∴ Mail classified as spam

Some of these features are more inclined towards spam than the others.

We can assign weight to each features depending on the priority. Weight vary from -1 to 1.

Features	Weight	Weight (Spam)	Weight (Not Spam)
F1: Look at who the email is addressed to	W1	-0.2	0.5
F2: Check the “from” email address	W2	-0.2	0.6
F3: Check the greeting	W3	-0.3	0.7
F4: Look at the Subject of the Email	W4	0.6	-0.3
F5: Email content are money, lottery, win, prizes	W5	0.7	-0.4
F6: Bad Grammar and Spelling	W6	0.5	-0.2
F7: No phone number in the signature	W7	-0.2	0.5

2 Moods - Happy, sad.

Today's mood - depends on yesterday's.

Based on her mood she wears diff. colors \leftarrow

Given the E & T_S tables.

(HMM)

W.i.T

H	R	G	B
S	0.8	0.1	0.1
S	0.2	0.3	0.5

T.i.T

H	S
H	0.7
S	0.5

Prob Day Moods vs today

We, as students won't go and ask tr. if she's happy or sad. \therefore These are HIDDEN LAYER.
but based on the color of Tshirt they're wearing, we can OBSERVE the hidden state.
Observations: Most likely moods

$$\begin{array}{l|l} \text{Day 1} & C_1 = G \\ & M_1 = ? \\ & M_2 = ? \\ & M_3 = ? \\ \text{Day 2} & C_2 = B \\ & M_1 = M_1, M_2 = M_2, M_3 = M_3 \\ \text{Day 3} & C_3 = R \\ & M_3 = ? \end{array}$$

$$\underset{\substack{\text{Max} \\ M_1, M_2, M_3}}{P}(C_1 = G, C_2 = B, C_3 = R \mid M_1 = M_1, M_2 = M_2, M_3 = M_3) \quad \text{all happening tog.}$$

Soln: $P(C_3 \mid C_2, C_1, M_3, M_2, M_1)$ But C_3 depends only on $M_3 \Rightarrow P(C_3 \mid M_3)$ find Prod. of these s.t. it's maximized
 $P(C_2 \mid M_2)$ $P(C_1 \mid M_1)$: Assume it gives S, S, H .
 $P(M_3 \mid M_2, M_1)$ $P(M_2 \mid M_1)$ $P(M_1)$

$$P(C_3 \mid M_3) = \frac{P(R \mid S)}{P(G \mid S) + P(B \mid S)}$$

$$0.5 + 0.5$$

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

0.5

Conditionally

Q. I went to the bank to deposit my money
 I sat by Riverbank & enjoyed the view
 The bank of river was covered in flowers.

word : Bank

Sense 1 : financial institution

Sense 2 : RiverBank

test : I deposited my money in bank & left to buy flowers

~~Actual Person1: 3-0-4~~

~~P(Bank2) = 0.4~~

$$P(\text{Bank1} \mid \text{Context}) = P("I", \text{"dep..."} \mid \text{Bank1}) + p(\text{Bank1})$$

$$P(\text{Bank2} \mid \text{Context2}) = P("I", \text{"dep..."} \mid \text{Bank2}) + p(\text{Bank2})$$

	Went	bank	deposit	money	River	Enjoyed View
Sent1	1	1		1	-1	
Sent2						
Sent3	1	1			1	0

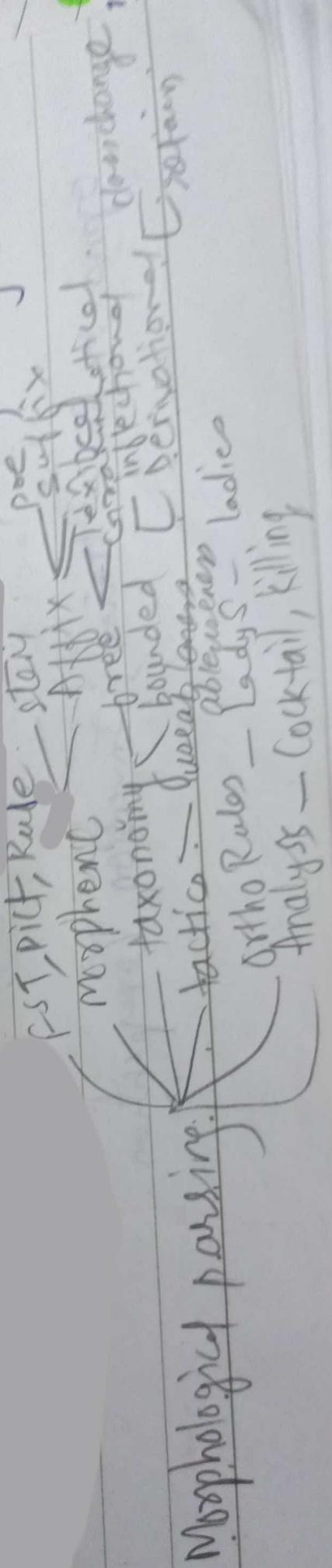
$$\begin{aligned}
 P(C1) &= P(\text{Context} = 1) + [P(\text{Went} \mid C=1) + P(\text{bank} \mid C=1) \dots] \\
 &= \frac{1}{3} * \left[\frac{2}{3} * \frac{1}{1} + \frac{1}{1} * 1 + \frac{1}{1} * 1 \dots \right]
 \end{aligned}$$

$$P(C2) = \dots$$

14.

Shift Reduce Parsing:

Stack	Input string	Action
\$.	The dog saw a man in the park \$	
\$ The	dog saw a man in the park \$	Det → "The"
\$ Det	dog saw a man in the park \$	
\$ Det dog \$ Det N	saw a man in the park \$ saw a man in the park \$	N → "dog".
\$ NP	saw a man in the park \$ a man in the park \$ a man in the park \$ men in the park \$ men in the park \$	NP → Det N .
\$ NP saw	a man in the park \$	V → "saw"
\$ NP V	a man in the park \$	det → "a".
\$ NP V a	men in the park \$	
\$ NP V det	men in the park \$	
\$ NP V det man	in the park \$	N → "man"
\$ NP V det N	in the park \$	NP → det N .
\$ NP V NP	in the park \$	
\$ NP V NP in	the park \$	P → "in"
\$ NP V NP p Det the	park \$	Det → "the"
\$ NP V NP Part Park	\$	N → "part"
\$ NP V NP P Det N	\$	NP → Det N
\$ NP V NP P NP	\$	PP → P NP .
\$ NP VP PP	\$	S → NP VP . P
\$ S		



3. Good Turing Algo.

- # $V = 9$
 - possible word pairs in bigram model (like in +let bigramTable count)
 $= 9 * 9 = 81$
 - seen distinct pairs = $7 + 2 + 1 = 10$.
- Total seen pairs = $0 \times 71 + 1 \times 7 + 2 \times 2 + 3 \times 1 = 14$
- # unseen pairs = $81 - 10 = 71$.

pair is
no. of times (b) repeating

$$\begin{matrix} C^+ \\ (1) \rightarrow NC \\ 1+7=7 \\ 71 \end{matrix}$$

C^-

N_C

$$71 \\ (81-7-2-1)$$

$$P^+ \\ C \text{ rank} \\ \# \text{ seen pairs} \\ \frac{0}{14} = 0$$

$$P^+ \\ C^+ \& N_C \\ \# \text{ seen pairs} \\ \frac{7}{14} = \frac{7}{71} * 14 =$$

$$2+2=4 \\ \frac{7}{7} \quad \frac{7}{7}$$

going abroad
abroad is
going to
study in
in the
field.

$$\frac{7}{14}$$

$$\frac{4}{14} \quad \frac{4+1+7}{7+14}$$

$$3+1=3 \\ \frac{2}{2}$$

2

2
is he
is going

$$\frac{4}{14}$$

$$\frac{3}{14} \quad \frac{3+1+2}{2+14}$$

3

1
He is

$$\frac{3}{14}$$

$$C^+ = (C+1)^+ \frac{N_{C+1}}{N_C}$$

or $P(\text{going}) = \frac{N_{C+1}}{N_C}$

$P(\text{is going}) \cdot \text{going abroad} \cdot \text{abroad to} \cdot \text{to study}$

$$= \left(\frac{3}{14} \div 2 \right) \cdot \left(\frac{4}{14} \div 7 \right) \cdot \left(\frac{7}{14} \div 71 \right) \cdot \left(\frac{4}{14} \div 7 \right) =$$

possible adjacent pairs \Rightarrow

he is/is ho/is going abroad abroad is/is going going to/to study study in in the field

3	2	1	1	2	1	1	1	1	1	1	1
---	---	---	---	---	---	---	---	---	---	---	---

Write all possible Bigrams down

he	is
is	he
is	going
going	abroad
abroad	is
is	going
to	study
study	in
in	the
the	field

→ 3 categories

- Count 1: 7 items
- Count 2: 2 items
- Count 3: 1 item

#words = 9, ∴ 9 * 9 = 81

#seen distinct pairs = 7 + 2 + 1 = 10

Total seen pairs = $7 \times 1 + 2 \times 2 + 3 \times 1 = 14$

#Unseen pairs = $81 - 10 = 71$

N_C, C^*, P, P^*

N_C Count of Bigrams	Size of Bigrams	$C^* = (C+1) * \frac{N_C+1}{N_C}$	$P = \frac{C^* - N_C}{\text{#seen pairs}}$	$\frac{P^*}{C^* - N_C}$ #seen pairs
71	0	$1 + \frac{7}{71} = \frac{7}{71}$	0	$\frac{7}{71}$
7	1	$2 + \frac{2}{7} = \frac{4}{7}$	$\frac{7+1}{14}$	$\frac{4}{7}$
2	2	$3 + \frac{1}{2} = \frac{3}{2}$	$\frac{4}{14}$	$\frac{3}{14}$
1	3	$4 + \frac{1}{1} = 5$	$\frac{3}{14}$	

$$P(\text{going}) = P(\text{is going}) \cdot P(\text{going abroad}) \cdot P(\text{abroad to}) \cdot P(\text{to study})$$

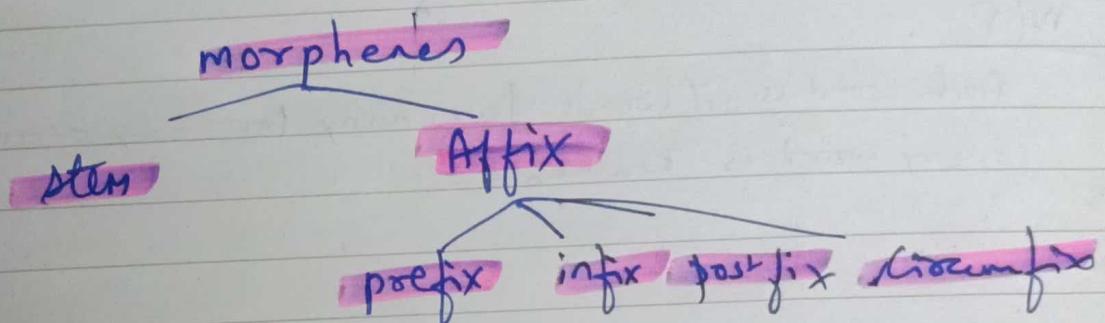
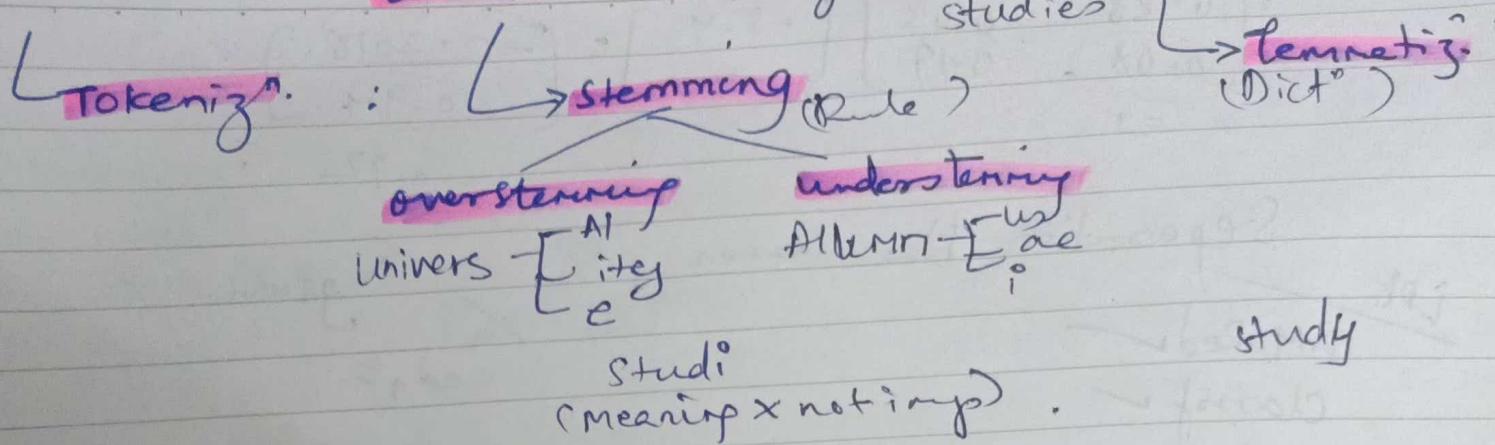
$$= \frac{1}{2}$$

$$\underline{P^* \div N_C}$$

25/10/23

Summary cheat sheet

word level Analysis



inflection
Derivation
compounding / Concatenating

Notebook
Laptop .

Regex

/ [abcd] /
[a-d]

[^ a]

carlt

[^ A - Z]

not a period

^ \ .

h, oh, - - .

Sa? j

o k h

y a +

ya, yaa, - - .

^ []
[] \$

anchor

morph. Models

dict
FSM DFA, NDFA, FST

SENTIMENT ANALYSIS

Stemming : get root stem. only 1 arg. the word stop

① Porter Stemmer

Eating → eat, history → histori, understanding → understand
Congratulations → Congratul, writes → write, fairly → fairli

② Lancaster Stemming

(choose?)

writing → writ
history → hist
finalized → fin
writes → writ

③ RegEx Stemmers

ing\$ ends w ing Remove kase

④ snowball Stemmer

(Best)

all ✓ history → histori, fairly → Fairli

App.
(CHATBOT)

Lemmatization : Get Root word

u provide 2 args . 1. word . 2nd. pos tag .

lemmatizer.lemmatize (word, pos = 'v')

programming, v ⇒ program
e) programming

LytK Algorithms

Membership Algorithm -

agar 1-4 Mei start symbol aqyng
 Then ✓ yes

$$S \rightarrow AB$$

$$A \rightarrow BB/A$$

$$B \rightarrow AB/B$$

check for $\overset{?}{abb}$.

	$A_{(1,1)}$	$S, B_{(1,2)}$	$A_{(2,2)}$	$S, B_{(1,4)}$
1				
2		$B_{(2,2)}$	$A_{(3,3)}$	$S, B_{(2,4)}$
3				
4			$B_{(3,3)}$	$A_{(3,4)}$
				$B_{(4,4)}$

$$\# \text{ cells} = \frac{n(n+1)}{2}$$

$$(1,2) = \underset{A}{(1,1)} + \underset{B}{(2,2)}$$

$$S \rightarrow AB \quad \therefore S, B$$

Bigrams -

$$(2,3) = \underset{B}{(2,2)} + \underset{B}{(3,3)}$$

$$B \rightarrow AB \quad \therefore A$$

$$(3,4) = \underset{B}{(3,3)} + \underset{B}{(4,4)} \quad \therefore A$$

$$(1,3) = \cancel{(1,1)} + \cancel{(3,3)}$$

$$(1,1) + (2,3) \rightarrow \underset{\phi}{\cancel{abb}}$$

Trigrams -

$$(1,4) = \cancel{(1,1)} + \cancel{(4,4)}$$

$$\underset{1234}{abb}$$

$$(2,4) = \cancel{(2,2)} + \cancel{(4,4)}$$

$$\underset{234}{bbb}$$

$$(1,2)(3) \underset{S, B}{\cancel{BB}} = \underset{\phi}{(SB)(BB)} = A$$

$$(2,3) \underset{A}{\cancel{A}} + (4) \underset{B}{\cancel{(4,4)}} = AB$$

$$S, B$$

$$(2) \underset{B}{\cancel{A}} (3,4) \underset{A}{\cancel{A}} = BA \times$$

$$(1,4) \Rightarrow \underset{1}{\overset{a}{\cancel{abb}}} \underset{234}{\cancel{bb}}$$

$$A \underset{S, B}{\cancel{A}} \Rightarrow (A, S) (A, B) \rightarrow SB$$

$$1,2 \underset{3,4}{\cancel{3,4}}$$

$$S, B \underset{A}{\cancel{A}} \Rightarrow (S, A) (B, A) \times \phi$$

$$(1,3) \text{ or } (1,2,3) \underset{4}{\cancel{4}}$$

$$A \underset{B}{\cancel{B}} \Rightarrow AB$$

$$S, B$$

Hey Pranav, can you add auto play
next video feature on codebasics.io



[12 1087 7 1870 1 2 ... 9]

Label
Encoding



label encoding if you have done

1 2 3 4 5
add, auto, and, assistance, are,
6 7 8 9
business, can, claim, codebsics.io,
10 11 12
data, 55, hey,
...
urgent
2000

VOCABULARY



	add	auto	and	assistance	are	business	can	...	urgent
Hey	[0	0	0	0	0	0	...	0]
Pranav	[0	0	0	0	0	0	...	0]
can	[0	0	0	0	0	0	1	0]
you	[0	0	0	0	0	0	0	0]
add	[1	0	0	0	0	0	0	0]
auto	[0	1	0	0	0	0	0	0]
...	...								



DMF

See here



	add, auto, and, assistance, need, I, immediately, help, urgent									
I	[0	0	0	0	0	1	0	0	0]
need	[0	0	0	0	1	0	0	0]	
help	[0	0	0	0	0	0	0	1]	
I	[0	0	0	0	0	1	0	0]	
need	[0	0	0	0	1	0	0	0]	
assistance	[0	0	0	1	0	0	0	0]	

Similar words do not have similar representation



the disadvantage number one.



if

too many words in Email .

Consumes too much memory & compute resources



So it consumes not only too much memory,

Visit codebasics.io for my practical, affordable video courses





Out Of Vocabulary (OOV) problem

numeric representation

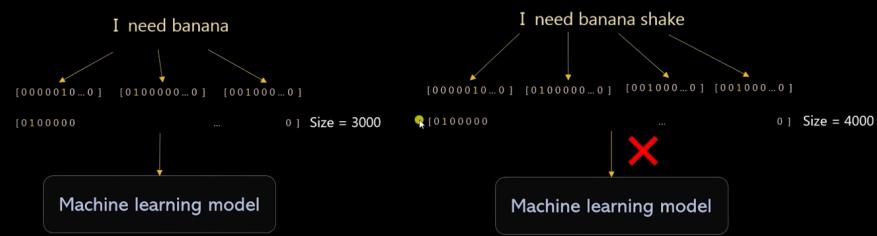


Out Of Vocabulary (OOV) problem



which can even represent the unknown
All Unknowns will have Same Numeric Representation





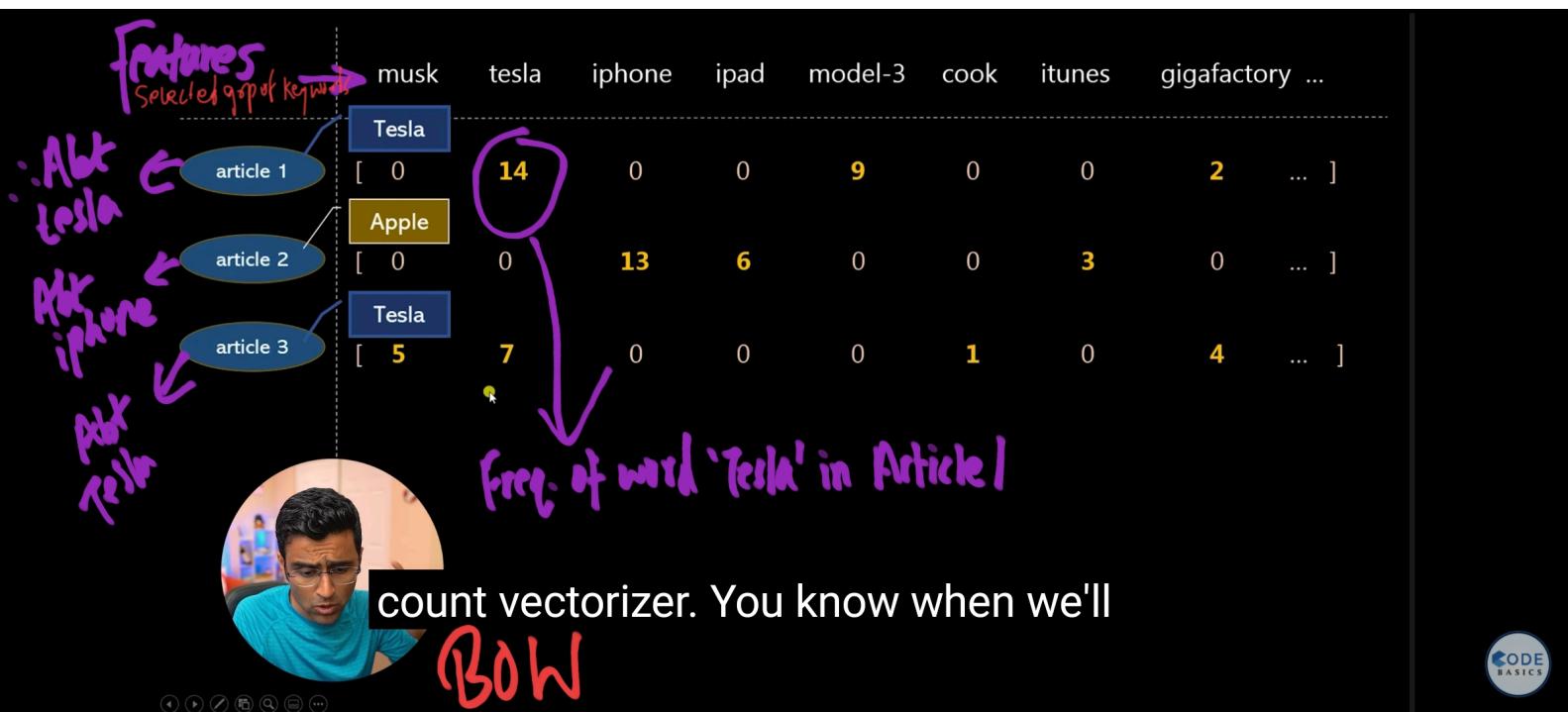
No fixed length representation



and that's the disadvantage number four.

All emails have diff. Sizes. A NN has a fixed size i/p.





	claim	free	help	health	hurry	so	play	money	up	volleyball
--	-------	------	------	--------	-------	----	------	-------	----	------------

Doc 1	play volleyball	[0	0	0	0	0	0	1	0	0	1]
-------	-----------------	---	---	---	---	---	---	---	----------	---	---	----------	---

Doc 2	Claim FREE money. Hurry up!	[1	1	0	0	1	0	0	1	1	0]
-------	-----------------------------	---	----------	----------	---	---	----------	---	---	----------	----------	---	---

Doc 3	Playing helps health so play!	[0	0	1	1	0	1	2	0	0	0]
-------	-------------------------------	---	---	---	----------	----------	---	----------	----------	---	---	---	---



text.

add, auto, and, assistance, need, I, immediately, help, ... urgent (100 k words)

I need help [0 0 0 0 1 1 0 1 ... 0] ← vector size=100k

Sparse Representation



That's called sparse presentation and
it may consume too much memory & compute resources



	add, auto, and, assistance, need, I, immediately, help, urgent											
I need help	[0	0	0	0	0	1	1	0	1	0]
I need assistance	[0	0	0	1	1	1	0	0	0	0]



emails is different. See here we have 1
Doesn't capture meaning or words properly



	add,	auto,	and,	assistance,	need,	I,	immediately,	help,	urgent	
I need help	[0	0	0	0	1	1	0	1	0]
I need assistance	[0	0	0	1	1	1	0	0	0]



numeric representation of both of these
Doesn't capture meaning or words properly

Dhaval sat on a sofa and
ate a samosa

Dhaval sat on a samosa
and ate a sofa



Meaning of a sentence is
determined by **order** of
words



Dhaval sat on a sofa and ate a samosa

BOW

Dhaval sat on a sofa and ate ...

bi-gram

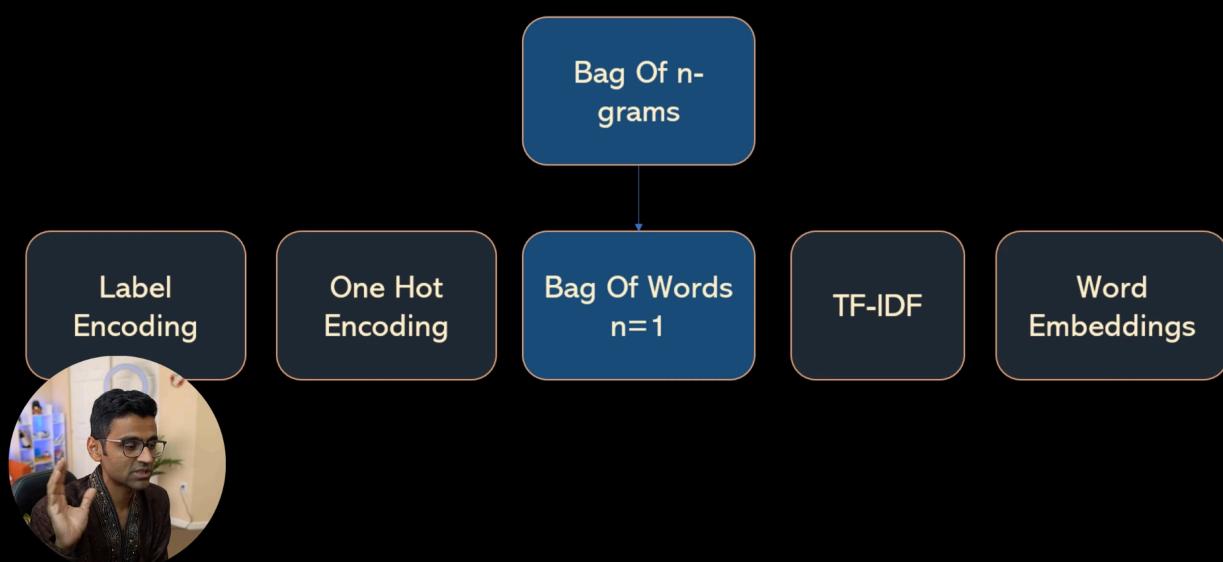
Dhaval sat sat on on a a sofa sofa and ...

tri-gram

Dhaval sat on sat on a on a sofa a sofa and ...



**There are various approaches of converting
text into vector**





Pre-processing

Thor ate pizza Doc 1
Thor eat pizza Doc 2
Loki is tall Doc 3
Loki tall
Loki is eating pizza
Loki eat pizza



	"Thor eat" "eat pizza" "Loki tall" "Loki eat"					
Doc 1	[1	1	0	0]
Doc 2	[0	0	1	0]
Doc 3	[0	1	0	1]



Word2vec

King – man + woman = Queen

We'll look all of this in detail, we'll



King

authority = 1

has tail = 0

rich = 1

gender = -1

[1,0,1,-1]

Horse

authority = 0

has tail = 1

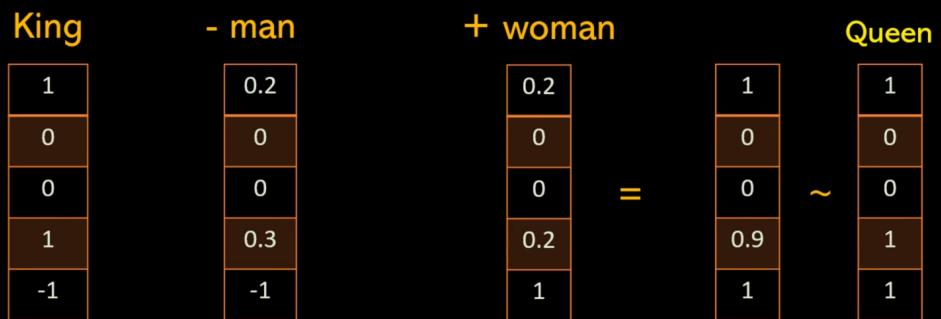
rich = 0

gender = -1

[0,1,0,-1]



	battle	horse	king	man	queen	..	woman
authority	0	0.01	1	0.2	1	...	0.2
event	1	0	0	0	0	...	0
has tail?	0	1	0	0	0	...	0
rich	0	0.1	1	0.3	1	...	0.2
gender	0	1	-1	-1	1	...	1

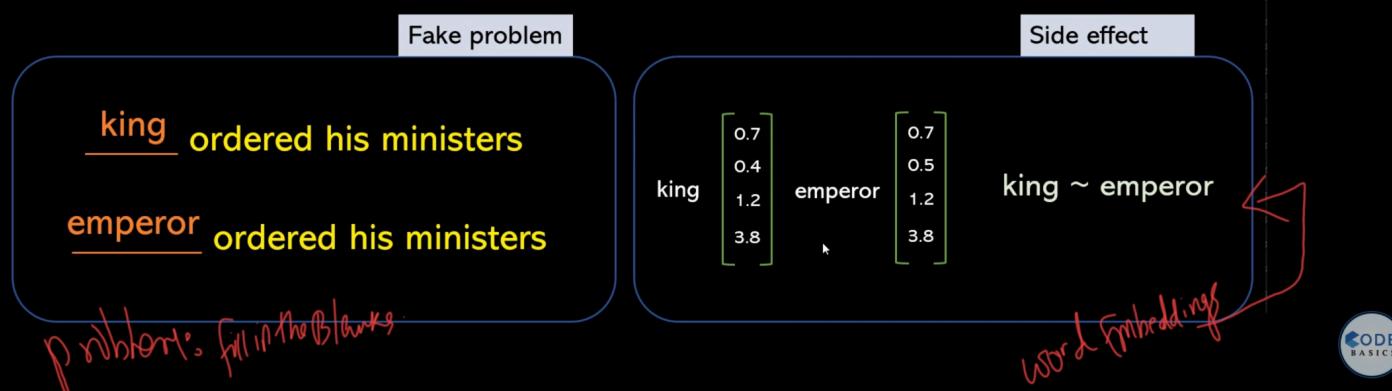


1. Take a **fake problem**
2. Solve it using neural network
3. You get **word embeddings** as a **side effect**

fake problem: fill in a missing word in a sentence



There lived a king called Ashoka in India. After Kalinga battle, he converted to Buddhism. This mighty king ordered his ministers to put together a peaceful treaty with their neighboring kingdoms. The emperor ordered his ministers to also build stupa, a monument with Buddha's teachings.



eating ____ is very healthy

table, angry, truck, **apple**, pizza, **walnut**

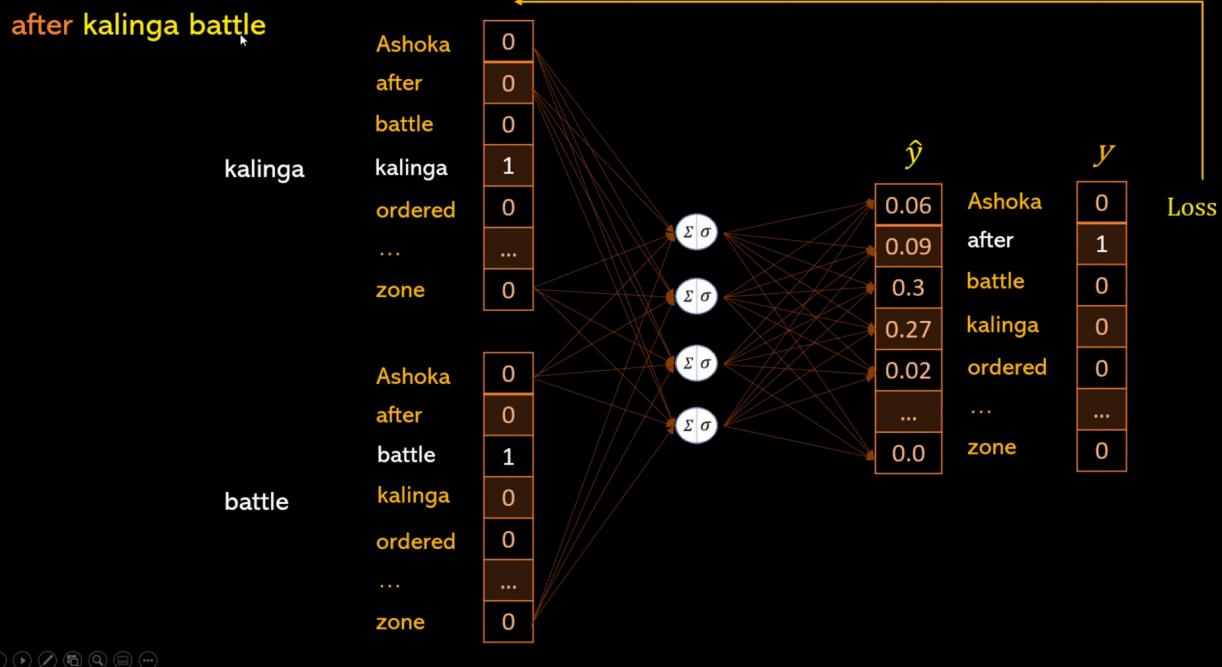


Meaning of word can be inferred by surrounding words

NASA launched ____ last month

table, angry, truck, **rocket**, apple, pizza







The first method, called **Continuous Bag of Words**, increases the context by using the surrounding words to predict what occurs in the middle.

Training Data

Troll 2 is great!
Gymkata is great!

