

NLP

→ an automated system capable of processing and analyzing text data

machine translation

→ conversion of one language to other

zipf's law

A relationship between the frequency of a word and its position in its list

$$f \propto \frac{1}{r}$$

(or)

$$f_r = k$$

The number of means  $m$  of a word obeys the law

$$m \propto \sqrt{f}$$

Heap's law:-

Let  $|V|$  size of vocabulary &  $N$  be no. of tokens

$$|V| = kN^\beta$$

Typically

$$k \approx 10 - 100$$

$$\beta = 0.4 - 0.6 \text{ (roughly sq. root)}$$

## Text pre-processing

Tokenization is process of separating a string of characters into words.

Sentence Tokenization  $\rightarrow$  splitting into sentences.

## Normalization

- Case folding
- stemming
- lemmatization

## Morphology:-

↳ internal structure of words

## Porter's algorithm

Py Enchant  $\rightarrow$  package (spelling correction)

- Insertion
- Deletion
- Substitution

M	b	5	6	5	4
U	5	4	5	4	5
I	4	3	4	5	6
D	3	2	3	4	5
E	2	1	2	3	4
M	1	2	3	4	5
#	0	1	2	3	4

MEDIUM

\*E\*NUM

d d s

N	9	8	9	10	11	12	11	10	9	8
O	8	7	8	9	10	11	10	9	8	9
I	7	6	7	8	9	10	9	8	9	10
T	6	5	6	7	8	9	8	9	10	11
N	5	4	5	6	7	8	7	8	9	10
E	4	3	4	5	6	7	6	7	8	9
T	3	4	5	6	7	8	7	8	9	10
N	2	3	4	5	6	7	6	7	8	9
I	1	2	3	4	5	6	5	6	7	8
#	0	1	2	3	4	5	4	5	6	7
#	E	X	E	C	U	T	I	O	N	

~~INTENTION~~  
~~EXECUTION~~

INTENTION  
+ EXECUTION

$O(mn) \rightarrow$  edit distance

$O(m+n) \rightarrow$  backtracking

dynamic edit distance

True Positive  
False +ve  
True -ve  
False -ve

Probabilistic Language models:-  
Machine translation  
Context based spell check  
Natural language generation.

$\langle s \rangle$  Helloworld  $\langle /s \rangle$

$N=1$

$$P(\langle s \rangle) \times P(\text{Hello} | \langle s \rangle) \times P(\text{world} | \langle s \rangle) \times P(\langle /s \rangle)$$

$N=2$

$$P(\langle s \rangle) \times P(\text{Hello} | \langle s \rangle) \times P(\text{world} | \text{Hello}) \times P(\langle /s \rangle)$$

$N=3$

$$P(\langle s \rangle) \times P(\text{Hello} | \langle s \rangle) \times P(\text{world} | \text{Hello}) \times P(\langle /s \rangle)$$

birrow

the dog smelled like a skunk

the dog  
dog smelled  
smelled like  
like a  
a skunk.

the dog sneled  
dog smelled like  
smelled like a  
like a skunk

11542, 703, 6268, 503, 179, 296

$$P(I | \langle s \rangle) = \frac{3}{5}$$

$$P(\text{Sam} | \langle s \rangle) = \frac{1}{3}$$

$$P(\langle s \rangle | \text{Sam}) = 0$$

$$P(\text{am} | I) = \frac{2}{4}$$

$$P(\text{do} | I) = \frac{2}{4}$$

$$P(\text{Sam} | \text{like}) = \frac{2}{3}$$

$$P(\text{like} | I, \text{do}) = 0$$

$$P(\text{Sam} | \text{not like}) = \frac{2}{3}$$

$$P(\text{Sam} | \text{do not like}) = \frac{2}{3}$$

allegation	$\frac{3}{7}$	$\frac{3 \times 3}{7 \times 7} = \frac{6}{49}$	$\frac{6}{28} \times \frac{1}{4} = 1.5$
reports	$\frac{2}{7}$	$\frac{5}{28}$	$\frac{5}{28} \times \frac{1}{7} = 1.25$
claim	$\frac{1}{7}$	$\frac{4}{28}$	$\frac{4}{28} \times \frac{1}{7} = 1$
request	$\frac{1}{7}$	$\frac{4}{28}$	$\frac{4}{28} \times \frac{1}{7} = 1$
attack	$\frac{0}{7}$	$\frac{3}{28}$	$\frac{3}{28} \times \frac{1}{7} = 0.75$
man	$\frac{0}{7}$	$\frac{3}{28}$	$\frac{3}{28} \times \frac{1}{7} = 0.75$
outcome	$\frac{0}{7}$	$\frac{3}{28}$	$\frac{3}{28} \times \frac{1}{7} = 0.75$

"data sparsity"

Balkoff & Interpolation :-

$$c(w_i/w_{i-2}, w_{i-1}) = \lambda_1 (w_i/w_{i-2}, w_{i-1}) + \lambda_2 (w_i/w_{i-1}) + \lambda_3$$

$$\lambda = \lambda_1 + \lambda_2 + \lambda_3$$

$$c^* = \frac{(c+1) N_{c+1}}{N_c}$$

- 3 allegation
- 2 reports
- 1 claims
- 1 request
- 0 attack
- 0 man
- 0 outcome

Given that in the training corpora the <sup>in</sup>grams

derived the alleg

estimate the modified count  $c^*$  and the modified

probability of the unseen <sup>in</sup>grams. Also find  $P^*$

for existing data. Use good Turing estimate for calculation of the unseen <sup>in</sup>grams.  $N_0 = 3$   $N_1 = 2$   $N_2 = 1$   $N_3 = 1$

$$c_0^* = \frac{1 \cdot 2}{3} = \frac{2}{3}$$

$$c_1^* = \frac{2 \cdot 1}{2} = 1$$

$$c_2^* = \frac{3 \cdot 1}{1} = 3$$

$$P_0^* = \frac{c^*}{N}$$

$$= \frac{2}{3} \times \frac{1}{7} = \frac{2}{21}$$

$$P_1^* = 1 \times \frac{1}{7} = \frac{1}{7}$$

$$P_2^* = 3 \times \frac{1}{7} = \frac{3}{7}$$

$$P_3^*$$

clock goes bit

Time on clock

Time on clock is twelve

clock - 3 time - 2  
goes - 1 on - 2  
tick - 1 ~~clock~~  
is - 1  
twelve - 1

The unseen unigrams in the training corpora are value, welcome, wait, use, good, turn, estimate. Find the modified probability values of  $P^*(value)$ ,  $P^*(time)$ ,  $P^*(tick)$

$$N_0 = 3 \quad N_1 = 4 \quad N_2 = 2 \quad N_3 = 1$$

$$c_0^* = \frac{1 \cdot 4}{3} = \frac{4}{3}$$

$$c_1^* = \frac{2 \cdot 2}{4} = \frac{4}{4} = 1$$

$$c_2^* = \frac{3 \cdot 1}{2} = 1.5$$

$$c_3^* = -$$

$$P_0^* = \frac{4}{3} \times \frac{1}{14} = \frac{1}{12} \quad \frac{4}{42}$$

$$P_1^* = \frac{4}{4} \times \frac{1}{14} = \frac{3}{37} \quad \frac{6}{101}$$

$$P_2^* = \frac{3}{2} \times \frac{1}{14} = \frac{1}{16} \quad \frac{3}{48}$$

$$P^*(value) = \frac{4}{33}$$

$$P^*(time) = \frac{3}{22}$$

$$P^*(tick) = \frac{1}{11}$$



If size of vocabulary  $V$  is given as input then no. of unigrams should be calculated as  $V$  (unigrams) =  $V$  (unigrams) +  $V$  (bigrams) +  $V$  (trigrams) + ...

$V^2$  (bigrams)  
 $V^3$  (trigrams)

Written Bell

$$P^* = \frac{C^*}{N} = \frac{1}{N+T} \left( \frac{T}{N+T} \right)$$

$$C^* = \frac{N}{N+T} \quad C=0$$

$$C^* = C \left( \frac{N}{N+T} \right) \quad C \neq 0$$

POSTagging

CRF - Conditional Random Field  
HMM - Hidden Markov Model.

denied the allegations = 3  
denied the claim = 1  
reports = 2  
request = 1

$T=4$

denied the attack  
denied the attack

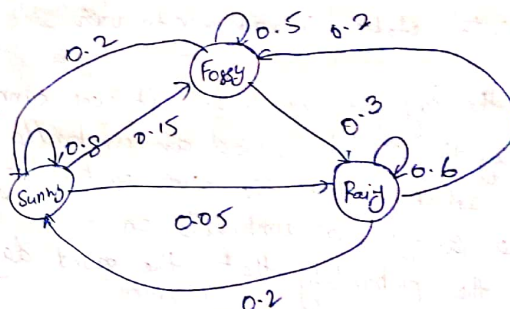
C	Nc	C*
0	3	$\frac{7}{3} \times \frac{4}{11} = \frac{28}{33}$
1	2	$\frac{2}{11} \times \frac{7}{11} = \frac{14}{121}$
2	1	$\frac{2}{11} \times \frac{7}{11} = \frac{14}{121}$
3	1	$\frac{2}{11} \times \frac{7}{11} = \frac{14}{121}$

POSTagging

classification sequence labels

Assume transitions in weather are indicated by a series of states which indicate what is the likelihood tomorrow's weather given today's.

	Sunny	Rain	Foggy
Sunny	0.8	0.05	0.15
Rain	0.2	0.6	0.2
Foggy	0.2	0.3	0.5



Given that it is sunny today what is the probability that tomorrow is sunny and day after is rainy.

$$P(w_2 = \text{sunny}, w_3 = \text{rain} / w_1 = \text{sunny}) = P(w_2 = \text{sunny} / w_1 = \text{sunny}) \times P(w_3 = \text{rain} / w_2 = \text{sunny})$$

$$= 0.08 \times 0.05 = 0.04$$

$$P(w_2 = \text{Rain} / \text{Foggy}) \times P(\text{rainy} / \text{rainy})$$

$$0.3 \times 0.6$$

$$= 0.18$$

$$P(w_3 = \text{rainy} / w_1 = \text{Foggy})$$

• (45)

$$0.2 \times 0.5 + 0.3 \times 0.6 + 0.5 \times 0.3$$

$$= 0.34$$

Hidden Markov Model is a state where one of the intermediate states is unknown to user

Suppose the day you were locked it was sunny the next day the caretaker carried an umbrella into the room assumes that the prior probability of the caretaker carrying an umbrella on any day is 0.5. What is the probability that the second day was rainy?

	Probability of carrying umbrella
Sunny	0.1
Rainy	0.8
Foggy	0.3

Markov

$$\frac{P(w_1, w_2, \dots, w_n, u_1, u_2, \dots, u_n)}{P(u_1, u_2, \dots, u_n)} = P(u_1, u_2, \dots, u_n / w_1, w_2, \dots, w_n) \times P(w_1, w_2, \dots, w_n)$$

$$P(u_1, u_2, \dots, u_n / w_1, w_2, \dots, w_n) = \prod P(u_i / w_i)$$

$$P(w_2 = \text{Rainy} / w_1 = \text{Sunny}, u_2 = \text{True})$$

$$= \frac{P(w_1 = \text{Sunny} / u_2 = \text{True})}{P(u_2 = \text{True} / w_1 = \text{Sunny})}$$

$$= \frac{P(w_2 = \text{Rainy} / w_1 = \text{Sunny}, u_2 = \text{True})}{P(w_1 = \text{Sunny})}$$

$$= \frac{P(u_2 = \text{True} / w_2 = \text{Rainy}) \times P(w_2 = \text{Rainy} / w_1 = \text{Sunny})}{P(u_2 = \text{True}) \times P(w_1 = \text{Sunny})}$$

$$= \frac{0.8 \times 0.05}{0.5} = 0.08$$

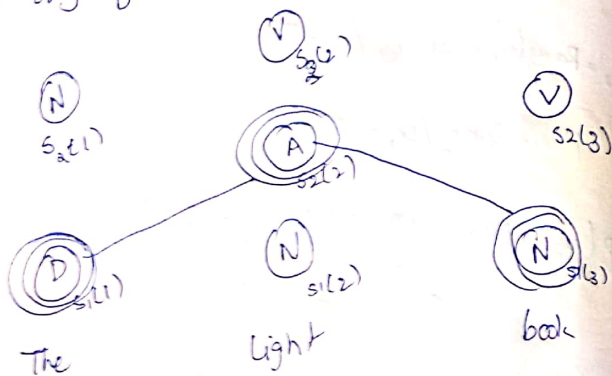
	Det	Noun	Adj	Verb
Det	0	0.5	0.3	0.0001
Noun	0	0.2	0.002	0.3
Adj	0	0.2	0	0.001
Verb	0	0.3	0	0.1

Tag transition probability  
 $P(t_i / t_{i-1})$

	D	N	A	V
The	0.3	0.1	0	0
light	0	0.05	0.2	0.06
book	0	0.03	0	0.001

Word Prob ( $w_i / t_i$ )

Given the phrase the light book. use HMM in context with Viterbi algorithm to find most probable tag sequence for the given phrase. The transition and the word emission probability and emission probabilities learnt from the training data are given in tables. Assume that the probability of any tag starting the sentence is same.



Level 1

$$s_1(1) = 0.25 \times 0.3 = 0.075$$

$$s_2(1) = 0.25 \times 0.1 = 0.025$$

$$s_{12} = \begin{cases} s_{11} \times P(N|D) \times P(\text{light}|D) \\ = 0.075 \times 0.5 \times 0.003 = 0.1125 \times 10^{-3} \\ s_{21} \times P(N|N) \times P(\text{light}|N) \\ = 0.025 \times 0.2 \times 0.003 = 1.5 \times 10^{-5} \end{cases}$$

$$s_{22} = \begin{cases} s_{11} \times P(A|D) \times P(\text{book}|A) \\ = 0.075 \times 0.3 \times 0.2 = 4.5 \times 10^{-3} \\ s_{21} \times P(A|N) \times P(\text{book}|A) \\ = 0.025 \times 0.002 \times 0.2 = 1 \times 10^{-5} \end{cases}$$

$$s_{23} = \begin{cases} s_{11} \times P(V|D) \times P(\text{light}|V) \\ = 0.075 \times 0.0001 \times 0.06 = 4.5 \times 10^{-8} \\ s_{21} \times P(V|N) \times P(\text{light}|V) \\ = 0.025 \times 0.3 \times 0.06 = 4.5 \times 10^{-4} \end{cases}$$

$$s_{31} = \begin{cases} s_{12} \times P(N|N) \times P(\text{book}|N) \\ = 0.1125 \times 10^{-3} \times 0.2 \times 0.03 = 6.75 \times 10^{-7} \end{cases}$$

$$s_{22} \times P(N|A) \times P(\text{book}|N) \\ = 4.5 \times 10^{-3} \times 0.2 \times 0.03 = 2.7 \times 10^{-5}$$

$$s_{23} \times P(N|V) \times P(\text{book}|N) \\ = 4.5 \times 10^{-4} \times 0.3 \times 0.03 = 4.05 \times 10^{-6}$$

$$s_{32} = \begin{cases} s_{12} \times P(V|N) \times P(\text{book}|V) \\ = 0.1125 \times 10^{-3} \times 0.3 \times 0.001 = 3.375 \times 10^{-8} \\ s_{22} \times P(V|A) \times P(\text{book}|V) \\ = 4.5 \times 10^{-3} \times 0.001 \times 0.001 = 4.5 \times 10^{-9} \\ s_{23} \times P(V|V) \times P(\text{book}|V) \\ = 4.5 \times 10^{-4} \times 0.1 \times 0.001 = 4.5 \times 10^{-8} \end{cases}$$



stemming  
morpheme

affixes → suffix / prefix / infix / circumfix  
 ↓  
 add at end to change meaning  
 at beginning  
 in between  
 add both sides

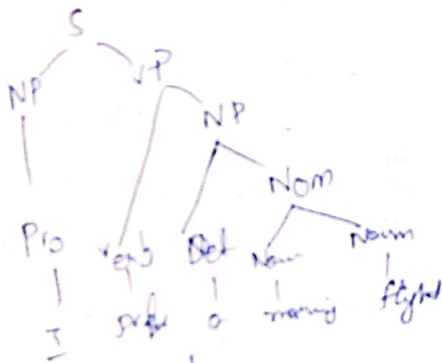
concatenative morphology (possible to break to get root)  
 non-concatenative morphology (not possible to break to get root)  
 e.g. - better

used: context based spellcheck, information extraction

inflectional morphology

derivational morphology

Context Free Language for English



Bracket Notation of parse tree

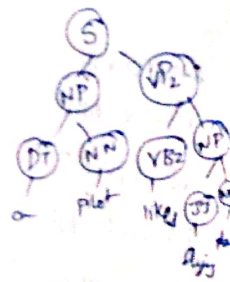
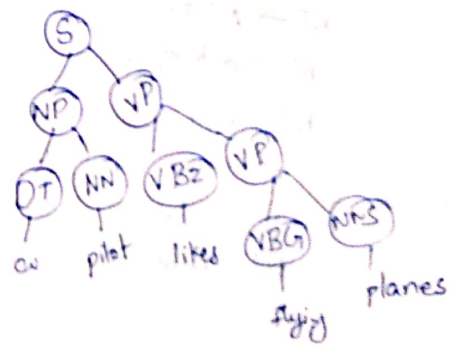
$[_S [NP [Pro I]] [VP [Verb Saw] [NP [Det a] [Noun man fly]]]]$

Chomsky Normal Form

Either, exactly two non-terminals on RHS  
 Or 1 terminal symbol on RHS

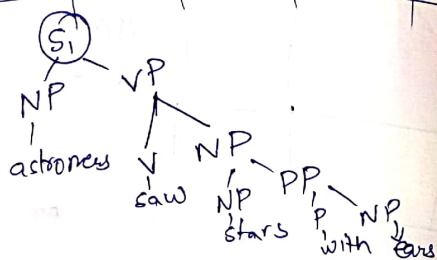
CKV Algorithm

	0	1	2	3	4	5	
	cu	pilot	likes	fly	planes		
DT		NP	-	-	S		S → NP VP VP → VBZ NP
		NN	-	-	-		VP → VBZ NP NP → DT NN NP → JJ NN
			VBZ	-	VP		DT → cu NN → pilot
				VBG	NP		VBZ → likes VBG → fly
					NNS		JJ → fly NNS → planes





astronomer's	saw	stars	with	ears
NP	-	S	-	S <sub>1</sub>
-	✓ NP	VP	-	VP V NP
-	-	N.P	-	NP
-	-	-	PP	PP
-	-	-	-	NP



S → NP VP  
 VP → V NP  
 VP → VP PP  
 PP → P NP  
 P → with  
 V → saw

NP → NP PP  
 NP → astronomer's  
 NP → ears  
 NP → saw  
 NP → stars  
 NP → telescope

Homonymy < Homophony  
Homographs

Polysemy

Synonymy

Antonyms

Hyponymy

Metonymy