

Why is NLP hard?

i) Lexical Ambiguity:

- ex: a) Will will will will's will?
- b) Real nose to put nose nose on her nose of noses.
- c) Buffalo buffalo Buffalo buffalo buffalo
buffalo Buffalo buffalo

ii) Language Ambiguity: Structural → diff interpretations of the same sentence

- ex: a) The man saw the boy with the binoculars
- b) Flying planes can be dangerous
- c) Hole found in the room wall; police are looking into it.

iii) Language imprecision & vagueness

- ex: a) It is very warm here [what is considered warm in China is diff from what is considered warm in US]

- b) Did your mom call your aunt last night?
↳ I'm still she must have. [basically saying idk]

Catalan Number: as the no. of phrases in a sentence ↑
the no. of interpretations in which they can be connected in a sentence ↑

Why is Language Ambiguous?

- goal in the production & comprehension of natural lang efficient communication
- allowing reusable ambiguity
 - permits shorter linguistic expressions
 - avoids language being overly complex
- language relies on ppl's ability to use their knowledge & inference abilities to properly resolve ambiguities.

Natural Language vs Computer Language

1. Cr doesn't have any ambiguity. They mean only one thing.
2. Formal programming lang are designed to be unambiguous
 - ↳ can be defined by a grammar that produces a unique parse for each sentence in the lang.
3. Programming languages are also designed for efficiency (deterministic) parsing.

Why else is NLP Hard?

→ idioms

- dark horse
- ball in your court
- burn the midnight oil

→ neologisms (used in social media as wbs)

- unfriend
- retweet
- Google/Skype/photoshop

Empirical Laws

- Function words: have little lexical meaning but serve as important elements in the structure of sentences
- to make the sentence more grammatical
- ex: words like 'a', 'the', pronouns etc
- Content words are various nouns & verbs that convey what are the important concepts in the sentence.
- mainly used for determining the structure of a sentence.

Function words are classless words
→ there are very specific grammatical categories there called function words
→ they are class-less words in the sense that we don't have never & never function words coming.
ex: prepositions, pronouns, auxiliary verbs, conjunctions, grammatical articles, particles etc, determiners.

Type-Token Distinction:

Type-token distinction is a distinction that separates a concept from the objects which are particular instances of the concept.

concept → Type
instance of } → Tokens
concept

ex: Will will

Type = 1	→	unique words
Tokens = 2	→	number of words (not unique)

Type-Token Ratio:

- TTR = $\frac{\text{Type}}{\text{tokens}}$ (no. of diff words)
tokens (no. of running words)
in a text or corpus
- this index indicates how often, on average, a new 'word form' appears in the text or corpus.
- High TTR → lots of new words coming
Low TTR → words in the corpus keep getting repeated

TTR by itself is not a valid measure of text complexity
→ the value varies with the size of text
→ for a valid measure, a running avg is computed on consecutive 1000-word chunks by text.

Zipf's Law:

- Count the freq of each word type;
- List the word types in decreasing

• Zipf's law gives the rel b/w the & its position in the list (its rank)

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

$$f \cdot n = \text{constant}$$

i.e., the 50th most common word is 3 times the frequency of the 150th

most common/rank	freq
50	1
150	1/3

$$f_1 = 3f_2$$

Another way to denote this,

Let P_n = denote the probability of N = denotes the total no. of

$$P_n = \frac{f_n}{N} \quad (\text{by definition})$$

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

$$\Rightarrow \frac{f}{N} = \frac{A}{N} \quad A = N \quad f \cdot n = A \cdot N \quad f \cdot n = K \quad A \approx 0.1$$

Zipf's Other Laws:

Correlation = Number of meanings
The no. of meanings m of a word

$$m \propto \sqrt{n}$$

given the first law,

$$m \propto \frac{1}{\sqrt{n}}$$

Combining

Correlation = Word length & word frequency is inversely proportional to their length

Zipf's Law:

- Count the freq of each word type in a large corpus
- List the word types in decreasing order of their freq
- Zipf's Law gives the relⁿ b/w the freq of a word (f) & its position in the list (its rank n) as:

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

$f \cdot n = \text{constant}$

i.e., the 50th most common word should occur with 3 times the frequency of the 150th most common word.

most common/rank	freq
50	f_1
150	f_2

$$f_1 = 3f_2$$

Another way to denote this,

Let P_n = denote the probability of word of rank n
 N = denote the total no. of word occurrences

$$P_n = \frac{f}{N} \quad (\text{by definition})$$

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

$$\Rightarrow \frac{f}{N} = \frac{A}{n}$$

$$\begin{aligned} f \cdot n &= A \cdot N \\ f \cdot n &= K \end{aligned} \quad \left. \begin{aligned} K &= A \cdot N \\ A &\approx 0.1 \end{aligned} \right\}$$

Zipf's Other Laws:

Correlation = Number of meanings & word freq
The no. of meanings m of a word obey the law:

$$m \propto \sqrt{n}$$

} Sage Law

given the first law,
 $m \propto \frac{1}{\sqrt{n}}$

Combining

∴

Correlation = Word length & word Frequency
word frequency is inversely proportional to
their length

Knaps Law:

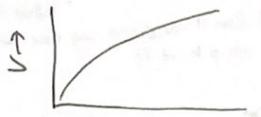
- gives relationship b/w size of vocab V and no. of tokens N

$$|V| = KN^\beta$$

$$K \approx 10-100$$

$$\beta = 0.4-0.6$$

no. of unique tokens ↑ with corpus size



no. of unique tokens ↑ with corpus size

Text Processing : Basics

1) Tokenization : split string into char of words

2) Sentence Segmentation

- ↳ challenge ↳ where does sentence end?
 - ↳ use binary classifier
 - ↳ use decision tree

3) Decision Tree

- ↳ feature can be separated into classes
- ↳ issues

- A) how to tokenize "Finland's"
- B)
- C) San-Francisco → "What's"
- D) m.p.h → ?? → one token or 2 tokens

4) Normalization

- ↳ indexed text & query ~~vectors~~ must have the same form
 - terms

ex.: U.S.A & USA should be matched.

5) Lemmatization

- ↳ reduce inflections or variant forms to base form
 - am, are, is → be
 - Car, Cars, Car's, Cars' → car

Spelling Correction: Edit Distance

- data obtained is noisy (with spelling errors)

How to use Edit Distance for Spelling correction?

Situation: There is a sentence

"I am writing this email on behalf of..."

Spelling of behalf is wrong.

You list down the words closest to the misspelled word:

- behalf
- behave ... etc.

Now, you have to choose 1 word from the given list of candidate words.

Isolated word error correction: trying to correct a word that is incorrect but not trying to use the context around it.

∴ we want to correct "behalf" w/o seeing the words before & after it.

∴ we use only candidate keys & the misspelled word.

How do we find out which candidate word is closest to the misspelled word?

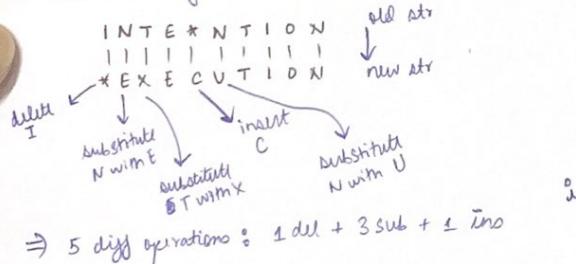
- using a distance metric

Edit distance: The simplest metric

→ minimum no. of edit operations needed to go from one string to another

→ various edit operations: insertion, deletion & substitution

ex: edit distance from 'intention' to 'execution'



⇒ 5 diff operations: 1 del + 3 sub + 1 ins

→ 2 diff variations to use:

(i) Each operation has a cost of 1 [Levenshtein]

$$\hookrightarrow \text{edit distance} = 5 \text{ op} \times \text{cost of 1} = 5$$

(ii) Substitution has a cost of 2

$$\hookrightarrow \text{edit distance} = (1 \text{ del} \times 1) + (1 \text{ ins} \times 1) + (3 \text{ sub} \times 2) \\ = 1 + 1 + 6 \\ = 8$$

How to find the minimum edit distance?

→ search for a path (sequence of edits) from the start $D(i-1, j-1) = 0$ to the final string

i) initial state: word we are transforming (intention)

ii) operators: ins, del, sub

iii) goal state: final string (execution)

iv) path cost: we want to minimize: the no. of edits

2 strings: X of length n

Y of length m

We define $D(i, j)$:

→ the edit distance b/w $X[1..i]$ and $Y[1..j]$
i.e., the first i changes X and the first j changes Y

The edit distn b/w X and Y is $D(n, m)$

NC

o o

o f

NO

i g

ii l

Na

we

we

re

ge

in

NI

w

si

w

pr

do

r

in

Er

ac

ar

an

an

an

an

an

an

Algorithm:

$$D(1, 0) = i$$

$$D(0, j) = j$$

For each $i = 1..M$

for each $j = 1..N$

$$D(i, j) = \min \begin{cases} D(i-1, j) + 1 & (\text{del}) \\ D(i, j-1) + 1 & (\text{ins}) \\ D(i-1, j-1) + 1 & (\text{sub}) \\ D(i-1, j-1) + 0; x(i) = y(j) & \end{cases}$$

Termination:

$D(n, m)$ is distance

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

$$\begin{aligned} D(i-1, j) &= 1 + 1 = 2 \\ D(i, j-1) &= 1 + 1 = 2 \\ D(i-1, j-1) &= 0 + 2 = 2 \\ D(i-1, j-1) &= 2 + 1 = 3 \\ D(i-1, j-1) &= 1 + 2 = 3 \\ D(i-1, j-1) &= 2 + 1 = 3 \end{aligned}$$

$$\text{ptr}(i, j) = \begin{cases} \text{LEFT} & (\text{ins}) \\ \text{DOWN} & (\text{del}) \\ \text{DIAG} & (\text{sub}) \end{cases}$$

Non-word spelling error detection

- any word not in a dict is an error
- the larger the dict the better

Non-word spelling error correction

- generate candidate: real words that are similar to the error word
- choose the best one
 - shortest weighted edit dist
 - highest noisy channel prob

Noisy Channel Model for Spelling Correction

we have an observation X , the misspelt word
we have the correct word W

now, we wanted to use W , but while going through a channel, we end up with a word X .

$$W \xrightarrow{\text{noisy channel}} X$$

Now, how do we find W , given X , the misspelt word.

$$\hat{W} = \underset{W \in V}{\operatorname{argmax}} P(W|X) = \underset{W \in V}{\operatorname{argmax}} P(X|W)P(W)$$

$$\hat{W} = \underset{W \in V}{\operatorname{argmax}} p(x|w)p(w)$$

words with } small edit distance
similar spelling } to error

words with } small edit dist of
similar pronunciation } pronunciation to error

Damerau-Levenshtein edit distance

minimum edit distance, where edits are:
ins, del, sub & transposition of 2 adj letters

Error	Candidate Correction	Correct Letters	Error Letters	Type	$P(X W)$
acres	acres	t	-	del	$\text{del} = 0.00017$
acres	erles	-	a	ins	$a \# = 0.00001$
acres	cares	ca	ac	trans	$aC = 0.000016$
acres	acees	c	n	sub	$r1C = 0.00000209$
acres	acres	o	e	sub	$e1o = 0.000093$
acres	ariso	-	s	ins	$o\$e = 0.0000321$
acres	acres	-	s	ins	$s\$e = 0.0000342$

$$P(X|W) = \begin{cases} \frac{\text{del}(w_{i-1}, w_i)}{\text{count}(w_{i-1}, w_i)} ; \text{del} \\ \frac{\text{ins}(w_{i-1}, w_i)}{\text{count}(w_{i-1}, w_i)} ; \text{ins} \\ \frac{\text{sub}(w_{i-1}, w_i)}{\text{count}(w_{i-1}, w_i)} ; \text{sub} \\ \frac{\text{trans}(w_{i-1}, w_i)}{\text{count}(w_{i-1}, w_i)} ; \text{trans} \end{cases}$$

N-Gram Language Model

- applications :
- Speech Recognition
 - Machine Translation
 - Context-Sensitive Spelling Correction
 - Natural Language Generation
 - Sentence Completion/Next word Prediction

Probabilistic Language Modelling

Goal: compute the prob of a sentence or seq of words

$$P(W) = P(w_1, w_2, w_3 \dots w_n)$$

related } : prob of an upcoming word

$$P(w_4 | w_1, w_2, w_3)$$

model that computes either of these is a language model

How to compute joint probability :

$P(\text{about, fifteen, minutes, from})$

$$P(B|A) = \frac{P(A, B)}{P(A)} \Rightarrow P(A, B) = P(B|A) \cdot P(A)$$

$$\Rightarrow P(A, B, C) = P(A) \cdot P(B|A) \cdot P(C|A, B)$$

$$P(A, B, C, D) = P(A) \cdot P(B|A) \cdot P(C|A, B) \cdot P(D|A, B, C)$$

$$P(\text{"about fifteen"}) = P(\text{about}) \cdot P(\text{fifteen}|\text{about}) \cdot P(\text{minutes}|\text{about fifteen}) \cdot P(\text{from}|\text{about fifteen minutes})$$

$$P(\text{after} | \text{about fifteen minutes from}) = \frac{\text{Count (about fifteen min from after)}}{\text{Count (about fifteen min from)}}$$

simplify this assumption:

$$P(\text{after} | \text{about fifteen minutes from}) \approx P(\text{after} | \text{from})$$

↓ current word ↓ prev word to current word

k^{th} order markov model:

$$P(w_1, w_2, \dots, w_n) = \prod_i P(w_i | w_{i-k}, \dots, w_{i-1})$$

$$\Rightarrow P(w_i | w_1, w_2, \dots, w_{i-1}) \approx P(w_i | w_{i-k}, \dots, w_{i-1})$$

→ An N-gram model uses only N-1 words as prior context

Unigram: $P(\text{apple})$

Bigram: $P(\text{apple} | \text{from})$

Trigram: $P(\text{apple} | \text{minutes from})$

An N-gram model is an $(N-1)$ order Markov model

Estimating N-gram probability

(i) Maximum Likelihood Estimate

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

Example:

<ss> I am here </s>

<ss> who am I </ss>

<ss> I would like to know </s>

$$P(I | \text{ss}) = 2/3$$

$$P(\text{am} | I) = 1/1$$

$$P(\text{here} | \text{am}) = 1/2$$

$$P(\text{who} | \text{ss}) = 1/1$$

$$P(\text{am} | \text{who}) = 1/1$$

$$P(I | \text{am}) = 1/2$$

$$P(\text{would} | I) = 1/3$$

⋮
and so on

Example 2:

Train sentences:

<ss> Three friends Amar, Akbar and

Antony are reading book </s>

<ss> Amar is reading Malgudi

day </ss>

<ss> Antony is reading a book by

Rahman </s>

<ss> Akbar is reading a detective

book </s>

count

Three	-	1	$\text{ss} = 4$
friends	-	1	$\text{ss} = 4$
Amar	-	2	
Akbar	-	2	
Antony	-	2	
and	-	1	
are	-	1	
reading	-	4	
book	-	2	
is	-	3	
Malgudi	-	1	
Day	-	1	
Detective	-	1	
by	-	1	
Rahman	-	1	

Test sentence

<ss> Amar is reading a book </s>

$P(\text{read} | \text{Amar}) \times P(\text{book} | \text{read})$

$\times P(\text{read} | \text{reading}) + P(\text{book} | \text{reading})$

$\times P(\text{reading} | \text{is}) + P(\text{book} | \text{is})$

$\times P(\text{is} | \text{reading}) + P(\text{book} | \text{is})$

$$\begin{aligned}
 &= \frac{\text{Count}(\text{Amer}, <\text{s}\rangle)}{\text{Count}(<\text{s}\rangle)} * \frac{\text{Count}(\text{is}, \text{Amer})}{\text{Count}(\text{Amer})} * \frac{\text{Count}(\text{reading}, \text{is})}{\text{Count}(\text{is})} * \\
 &\quad \frac{\text{Count}(\text{a}, \text{reading})}{\text{Count}(\text{reading})} * \frac{\text{Count}(\text{book}, \text{a})}{\text{Count}(\text{a})} * \frac{\text{Count}(</\text{s}\rangle, \text{book})}{\text{Count}(\text{book})} \\
 &= \frac{1}{4} * \frac{1}{2} * \frac{3}{8} * \frac{2}{4} * \frac{1}{2} * \frac{2}{3} \\
 &= \frac{1}{4} * \frac{1}{2} * \frac{2}{4} * \frac{1}{2} * \frac{2}{3} \\
 &= \frac{2}{24} = \frac{1}{12} = 0.083
 \end{aligned}$$

Practical issues

→ multiplication of many prob values
 ↳ do log as it avoids underflow & adding is faster than multiplying

$$\log(p_1 * p_2 * p_3) = \log p_1 + \log p_2 + \log p_3$$

→ handling zeros
 ↳ use smoothing

Evaluation of Language Models, Basic Smoothing

- One way to evaluate is by assigning high prob to real sentences (prob abd) over ungrammatical ones (rarely abd)
- but this has to be done manually
↳ can it be done automatically?
↳ by dividing dataset into training & testing

Extrinsic Evaluation:

- Learn 2 models A & B.
- use each model for one or more tasks such as spelling checker, speech recognizer etc
- compare performance of A & B
- whichever gives better per., that model is preferable (task-based evaluation)

Intrinsic Evaluation : Perplexity:

Shannon Game (based on Intuition)
↳ how well one can predict the next word

- we have 2 models.
- we apply them here & the models will give diff words to predict the next word in a sentence.
- whichever model predicts better is preferable
→ in this case, unigram models are not preferable
↳ don't use context around a word
- in this case, a better model will be one that assigns a high prob to the actual word.

Perplexity:

- we want a model that best predicts an unseen test data

Perplexity is the inverse prob of the test data, normalized by the no. of words.

→ Lower the perplexity, better the model

$$PP(w) = P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}}$$

where does the language model come into all this?

$$PP(w) = \left(\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})} \right)^{\frac{1}{N}} \quad [\text{applying chain rule}]$$

→ in the form of a unigram model.

so in we can see how many our context words we want to predict the next word.

Then, we can put in any bay model & fit it.

ex: For bigram model,

$$PP(w) = \left[\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})} \right]^{\frac{1}{N}}$$

ex: consider a sentence consisting of N random digits
The model assigns a prob. $P = \frac{1}{10}$ to each digit.
What's the perplexity of this sentence?

$$\begin{aligned} PP(w) &= P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}} \\ &= P\left(\frac{1}{10}, \frac{1}{10}, \dots, \frac{1}{10}\right)^{\frac{1}{N}} \\ &= \left(\frac{1}{10}\right)^N \end{aligned}$$

$$PP(w) = 10$$

Next choose a random bigram (w_i, x) as per its prob to continue the sent.
 $(w_1, w_1) \rightarrow [1]$ again choose a word based on this prob.
 $(w_1, w_2) \rightarrow [1]$
 $(w_2, w_3) \rightarrow [b]$ then again continue the process with
 $(w_2, w_4) \rightarrow [1]$ (next candidate)
(word, possibility)

Now, when do we stop? How to say wt the sentence is complete?
→ when you sample a word that is end of a sentence $(w_i, </s>)$

Now, in a scenario, (for example) if we take a corpus, we will find there are 844M bigrams. Out of them only 300K bigrams are real occurring.

This means that 300K bigrams have a while rest have prob = 0.

Now if you take a trained testing d. try to apply this model that was on the above corpus, then a bigra is present in testing data but not in data will get assigned a prob = 0.

If you try to find perplexity of t perplexity is multiplication of all diff multiplication o/p = 0.
We want the perplexity to have p greater than 20.

↓ How to tackle this issue of By Smoothing
↳ assigning diff prob va

ex: Training set:

- ... denied the allegations
- ... denied the reports
- ... denied the claims
- ... denied the requests

Due to testing data, not being $P(\text{offer} | \text{denied th}) = 0$

due to prob = 0, denominator \Rightarrow word can

$$\langle s \rangle, w_1 \rightarrow \left[\begin{array}{c} 1 \\ 0 \\ 0 \end{array} \right]$$

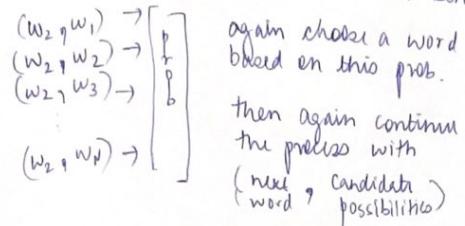
$$\langle s \rangle, w_2 \rightarrow \left[\begin{array}{c} 0 \\ 1 \\ 0 \end{array} \right]$$

$$\langle s \rangle, w_N \rightarrow \left[\begin{array}{c} 0 \\ 0 \\ 1 \end{array} \right]$$

choose w_i based on max prob

This is called Sampling a word from multinomial dist"

Next choose a random bigram (w_i, x) as per its prob to continue the sent.



Now, when do we stop? How to say when the sentence is complete?

→ when you sample a word that is the end of a sentence ($w_i, </s>$)

Now, in a scenario, (for example)

if we take a corpus, we will find that there are 844M bigrams. Out of all of them only 300K bigrams are real or freq occurring.

↳ This means that 300K bigrams have a prob > 0 while rest have prob = 0.

Now if you take a testing data & try to apply this model that was trained on the above corpus, then a bigram which is present in testing data but not in training data will get assigned a prob = 0.

If you try to find perplexity of testing data, perplexity is multiplication of all diff prob, then multiplication o/p = 0.

We want the perplexity to have preferably value greater than 200.

↓ How to tackle this issue of zeros

By Smoothing

↳ assigning diff prob values to get around zeros.

Ex: Training set:

- denied the allegations
- denied the reports
- denied the claims
- denied the requests

Due to testing data, not being part of training data

$$P(\text{offer} | \text{denied the}) = 0$$

due to prob = 0, denominator of perplexity = 0
 \Rightarrow and compute perplexity

Idea of Smoothing:

$$P(w| \text{denied the})$$

where w based on training & testing set.

	word count	prob
allegations	3	3/7
reports	2	2/7
claims	1	1/7
requests	1	1/7
	7	

what abt the prob of other bigram words like offer & loan?

→ First, assign them a prob of zeros

→ Second, idea of smoothing is can I

take some prob mass from each of the 4 training set words & assign them to the raw testing set words that can be said as unseen data.

so, if you remove 0.5 prob mass from the 4 words, you get a value of 2, which is prob 2/7, which you distribute among all other words in the corpus.

This dist' can be done in multiple ways

Testing set:

- denied the offer
- denied the loan

Laplace Smoothing

- also called "add-one estimation"
- here, pretend as if we saw every N-gram one more time than we actually did.
So, for those N-grams that you have not seen at all, you pretend you've seen them once
- ↳ just add 1 to all the counts

$$\text{MLE estimate for } P_{\text{MLE}}(\text{bigram}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

For Laplace smoothing this,
we add 1 to numerator (for each possible bigram)

$$\Rightarrow c(w_{i-1}, w_i) + 1$$

but to ensure that the prob sums upto 1, we have to modify the denominator such that they are normalized, so what to add to the denominators.

↳ how many diff bigrams will be there for which we will be adding 1?

$$\hookrightarrow V \text{ (no. of words in the vocab)}$$

so, we will add one 'V' times.

So, to normalize, we'll have to add a V to the denominators

$$\therefore \text{Add-one estimate} : P_{\text{add-1}}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$

Variand of Add-1 : Add-K :

$$P_{\text{add-K}}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) + K}{c(w_{i-1}) + KV}$$

NOTE:
if $KV = m$
then $K = m(\frac{1}{V})$

Unigram Prior Smoothing :

$$P_{\text{unigram prior}}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) + m p(w_i)}{c(w_{i-1}) + m}$$

Computational Morphology:

- Computational Morphology:

 - The internal structure of words, how words are built up from smaller meaning units are called morphemes.

ex: dogs = dog + S ← plural marker on nouns

Allomorphs:

- given a word, using morphemes, you can:
the word to its opposite
 - variants of the same morpheme, but can't be replaced by one another \Rightarrow means can't be used for ir-rational
 - ex: un-happy
 - ir- rational
 - in - comprehensible
 - im- possible

Bound Morphemes

- can't appear as a word by itself
 - -S (dog-s)
 - -ly (quick-ly)
 - -ed (walk-ed)

Free Morphemes:

- can appear as a word by itself
 - often can combine with other morphemes too
 - house -s
 - walk (walk-ed)
 - of
 - the
 - or

Stems & Affixes.

Stems (nests): the core meaning below units

Affixes: bits & pieces adhering to stems to change their meanings & grammatical functions

Types of Affines :

- ① Prefix: before word
ex: un-, anti-, pre-
 - ② Suffix: after word
ex: -ity, -ation, -ly, -ing
 - ③ Infix: in between word
ex: in Sanskrit, vi-n-de

Content Morphemes:

- carry some semantic content
Ex: con-, -able, un-

Functional Morphemes:

- provide grammatical info

ex: -s (plural),

Inflectional Morphology:

ex: Walk + ing → Walking

- makes changes in grammar such as number, tense, gender, case
 - creates new forms of same word → by

Derivational Morphology

Verb noun

- create new words by changing pos
ex: logic, logical, illogical, illogically

Applications :

- ① text to speech synthesis
 - ② search & information retrieval
 - ③ machine translation
 - ④ grammar correction

Intro to POS Tagging:

- given a sentence, identify the POS of each word

How many?

① Open class words (content words)

- nouns, verbs, adjectives, adverbs
- mostly content-bearing: they refer to objects, actions, & features in the world
- open class, since new words are added all the time.

② Closed class words

- pronouns, determiners, prepositions, connectives
- limited no. of these
- mostly functional: to tie the concepts of a sentence together.

POS examples

N	Noun	chair, bandwidth, pacing
V	Verb	study, debate, munch
Adj	Adjective	purple, tall
Adv	Adverb	unfortunately, slowly
P	Preposition	of, by, to
Prp	Propron	I, me, mine
Det	Determinant	the, a, that, these

→ to do POS tagging, a standard set needs to be chosen
 ↳ Universal Treebank tagset (as tag)

ex: The grand jury commented on a no. of other topics
 dt JJ NN VBD IN DET NN IN JJ NNS

why is POS tagging hard?

- words often have more than 1 POS

ex: The back door : JJ
 On my back : NN
 Win the voters back : RB
 Promised to back the bill : VB

POS Tagging: 2 Approaches

① Rule-Based Approach

- Assign each word in the input a list of potential POS tags
- Then winnow down this list to a single tag using hand-written rules

② Statistical Tagging

- Get a training corpus of tagged text, learn the transformation rules from the most frequent tags (TBL tagger)
- Probabilistic: Find the most likely sequence of tags T for a sequence of words W

Probabilistic Tagging

Problem: we have some data of $\{(d_i, c_i)\}$ of paired observations d and hidden classes c

Different instances of d and c

• Parts-of-Speech Tagging

→ words are observed and tags are hidden

• Text Classification

→ sentences/documents are observed & the category is hidden

→ categories can be positive/negative for sentiments

→ sport, politics, business for documents

Generative Models (Joint)

Generate the observed data from st stuff
 i.e., put a prob over the observations given the class
 ↳ $p(d, c)$ in terms of $p(d|c)$

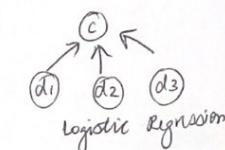
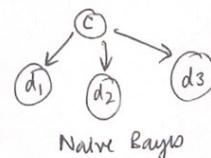
ex: Naive Bayes
 Hidden Markov Models

Discriminative (Conditional) Models

Take the data as given, and put a prob over hidden structure given the data

↳ $p(c|d)$

ex: logistic regression
 maximum entropy models
 conditional random fields



→ A joint model gives prob $p(d, c)$ and tries to maximize this joint likelihood

→ A conditional model gives prob $p(c|d)$ taking the data as given and modeling only the conditional prob of the class

(d, e) of
classes

Hidden Markov Model (for POS Tagging)

HMM is a probabilistic model to infer unobserved data / info from observed data

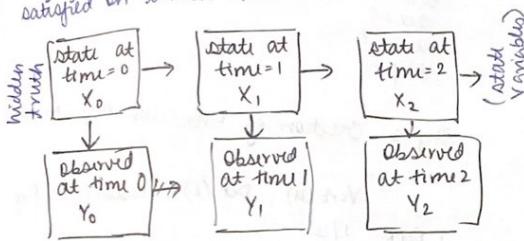
ex:
in mobile, sometimes you mistype the char
not to what you wanted to.
→ the char you mistyped is the observed data
→ one you intended to type is unobserved data

Observed
negative
for document

Similarly,
there are many cases where what is observed
is not the truth & the HMM is one way to
give us the hidden truth.

NOTE:
there are some assumptions that need to be
satisfied on which HMM works.

n st steps
etc



puta
in the

for each observed data, there's a hidden truth
(hidden truth = state variable)

then are a sequence of states (hidden truths)
linked with each other.

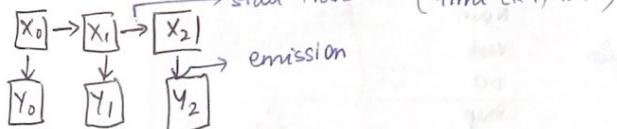
arrow = dependence
↳ any state depends ONLY on some of
its previous state(s).

a state that doesn't depend on any prev state] = independent

a state depends on its immediate prev state] = simple markov chain

a state depends on its 'n' prev states] = n-order markov chain

State Transition, Emission & Initial State



i) assume there are M possible states to choose from $\{1, 2, \dots, M\}$

ii) transition prob from state i to j is a_{ij}

iii) NOTE:

transition can happen from any one state to another out of M possible states
so there are $M \times M$ possibilities which can be arranged in a state transition matrix A

$$\text{From } \begin{bmatrix} a_{00} & a_{01} & a_{02} & \dots & a_{0m} \\ a_{10} & a_{11} & a_{12} & \dots & a_{1m} \\ \vdots & & & & \\ a_{m0} & a_{m1} & a_{m2} & \dots & a_{mm} \end{bmatrix} \rightarrow \text{All } a_{ij} \text{ b/w 0 & 1}$$

for a stationary markov chain,
matrix A is constant over time
else,

its diff for each timestep

$$a_{ij} = P(X_k=j | X_{k-1}=i)$$

$$S_K^T = S_{K-1}^T \cdot A \quad \text{where } S_k = \begin{bmatrix} P(X_k=1) \\ P(X_k=2) \\ \vdots \\ P(X_k=m) \end{bmatrix}$$

↓
multiplying the prob dist of time $(K-1)$ by A updates the prob dist to next time step K .

emission matrix B for storing the prob of a state;
resulting in an observed value j .
Observable is assumed to have k possible values, meaning it could take any one of $\{1, 2, \dots, k\}$

$$\text{From } \begin{bmatrix} b_{00} & b_{01} & \dots & b_{0k} \\ b_{10} & b_{11} & \dots & b_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ b_{m0} & b_{m1} & \dots & b_{mk} \end{bmatrix} \quad \text{To}$$

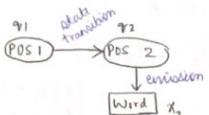
b_{ij} represents the prob of state i to emit observable j .

$$b_{ij} = P(Y_k=j | X_k=i)$$

$$O_k = S_k * B \quad \text{when } O_k = \begin{bmatrix} P(Y_k=1) \\ P(Y_k=2) \\ \vdots \\ P(Y_k=m) \end{bmatrix}$$

POS Tagging, HMM, Viterbi Alg

Generative Model: HMM



here, we want to determine the POS tag for word x_2

q_1 = prev state

q_2 = current state

we want to determine to which category the word emitted by POS 2 belongs to?

HMM is a comb of ST & SE

For ST, we use a bigram [can go for tri/n-gram too]
For SE prob, we are using the bayes rule.

Formula of HMM

$$P(\text{POS}_2 | (\text{word}, \text{POS}_1)) = P(\text{word} | \text{POS}_2) * P(\text{POS}_2 | \text{POS}_1)$$

← bayes rule

bigram
obtained from emission prob matrix

data is obtained from state transition matrix

So, we are determining the POS tag for x_2 based on POS_1

Emission Probability Matrix:

- From a given corpus, identify unique tokens.
Also identify total POS tag present in the corpus.
- Assign correct POS tag for every token
- Arrange tokens in horizontal & POS tags in vertical fashion. (We can change this order also)
- This will provide data for Bayes rule emission prob
- For N -tokens, M parts of speech tags, size of the emission matrix will be $N \times M$.

State Transition Prob:

- Use POS tags found from emission prob matrix
- Create matrix size of $M \times M$
- Use bigram concept & fill the entries in a table
- This will provide data for state trans prob

Need for Viterbi Decoding Alg:

- It's based on DP
- DP breakdowns the prob into sub-parts
- It stores the result for future purpose
- ∴ no need to compute the result again

Training Corpus with correct POS tag:

<S> Book a car </S>

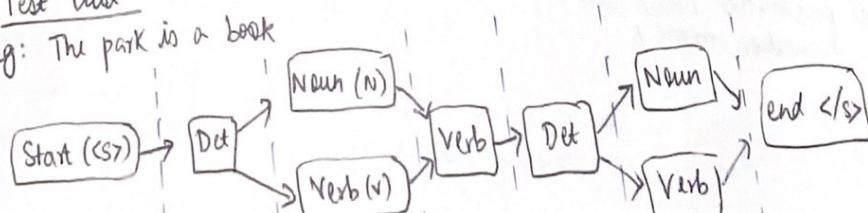
<S> park the car </S>

<S> The book is in the car </S>

<S> The car is in a park </S>

Test Data

e.g.: The park is a book



<S> The park is a book </S>
total no. of possibilities = $(1 \times 2 \times 1 \times 1 \times 2) = 4$

⇒ Viterbi alg is used to find the efficient path

Step 1: Assign correct POS tag to training corpus

<S>	Book	a	car	</S>	
<S>	Park	the	car	</S>	
<S>	The	book	is	the	car
<S>	The	car	is	in	a park
					noun

I) The First Word: "The"

$$\begin{aligned} q_1 \rightarrow \text{det} \\ \text{det} \rightarrow \text{The} \\ \text{The} \rightarrow x_2 \end{aligned}$$

$$P(q_1 | \text{S}) \cdot P(x_2 | q_1) + P(q_2 | q_1)$$

$$P(\text{det} | \text{the}, \text{S}) = P(\text{the}/\text{det}) + P(\text{det})$$

$$= \frac{1}{4} + \frac{2}{4} = \frac{1}{3}$$

emission prob
state trans

II) The Second Word: "Park"

previous prob = $\frac{1}{3}$

$$\begin{aligned} q_2 \rightarrow \text{noun} \\ \text{noun} \rightarrow \text{park} \\ \text{park} \rightarrow x_2 \end{aligned}$$

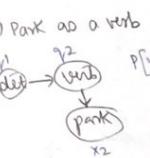
$$P(\text{noun} | \text{park}, \text{det}) = P(\text{park})$$

$$= \frac{1}{6}$$

Step 2: Creation of Emission Prob Matrix

	Verb (V)	Dt (D)	Noun (N)	Prep (P)
1. Book	1/4		1/6	
2. o		2/6		
3. Car			4/6	
4. Park	1/4		1/6	
5. the			4/6	
6. is	2/4			
7. in				2/2
Σ	4	6	6	2

$$P(\text{word} | \text{class}) = \frac{P(\text{class}, \text{word})}{P(\text{class})}$$



$$P(\text{verb} | \text{park}, \text{dt}) = P(\text{v})$$

$$= \frac{1}{4}$$

sometimes & path
2 states, value 1
we do something
∴ we multiply

$$\therefore P(\text{verb} | \text{park}, \text{dt}) = \frac{1}{4} + \frac{1}{6}$$

Now, compare $\frac{1}{1200}$ & $\frac{1}{72}$

$\frac{1}{72} > \frac{1}{1200}$ ∴ park is current state

Step 3: State Trans Prob Matrix

	Noun	Verb	Dt	Prep	Σ	total no. of word count
<S>	2/4	2/4			4	
Noun	2/6				6	
Verb		2/4	2/4		4	
Dt	1/6				6	
Prep			2/2		2	

verb follows
<S> twice

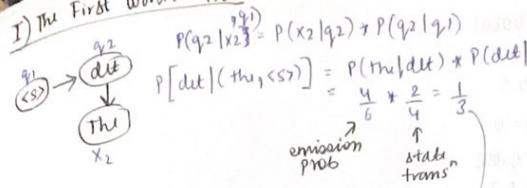
<S> occurs
after verb
4 times

III) The
prev prob
 q_1
noun →

IV) The four
prev prob.
 q_1
verb →

$$P(w_i | w_{i-1}) = \frac{\text{Count}(w_{i-1}, w_i)}{\text{Count}(w_{i-1})}$$

I) The First Word: "The"



II) The Second Word: "Park"

previous prob = $\frac{1}{3}$

i) Park as a noun

$$P[\text{noun} | (\text{park}, \text{det})] = P(\text{park} | \text{noun}) + P(\text{noun} | \text{det})$$

$$= \frac{1}{6} * \frac{1}{6} = \frac{1}{36}$$

$$P[\text{noun} | (\text{park}, \text{det})] = P(\text{noun} | \text{park}, \text{det}) * \text{prev prob}$$

$$= \frac{1}{6} * \frac{1}{3} = \frac{1}{18}$$

ii) Park as a verb

$$P[\text{verb} | (\text{park}, \text{det})] = P(\text{park} | \text{verb}) + P(\text{verb} | \text{det}) * \text{prev prob}$$

$$= \frac{1}{4} * 0 * \frac{1}{3} = 0$$

sometimes a path is efficient but in b/w 2 states, value becomes zero. So that time we do smoothing.

∴ we multiply by $(1/100)$

$$\therefore P(\text{verb} | \text{park}, \text{det}) = \frac{1}{4} * \frac{1}{100} * \frac{1}{3} = \frac{1}{1200}$$

Now, compare $\frac{1}{1200} > \frac{1}{18}$

$\frac{1}{18} > \frac{1}{1200}$ } ∵ park is a noun
noun > verb current prob = $1/18$
 for next step, it becomes prev prob

III) The third word: "is"

prev prob = $1/18$

$$P[\text{verb} | (\text{is}, \text{noun})] = P(\text{is} | \text{verb}) * P(\text{verb} | \text{noun}) * \text{prev prob}$$

$$= \frac{2}{4} * \frac{2}{6} * \frac{1}{18} = \frac{1}{108}$$

IV) The fourth word: "a"

prev prob = $1/108$

$$P[\text{det} | (\text{a}, \text{verb})] = P(\text{a} | \text{det}) * P(\text{det} | \text{verb}) * \text{prev prob}$$

$$= \frac{2}{6} * \frac{2}{4} * \frac{1}{108} = \frac{1}{648}$$

V) The Fifth word: "book"

prev prob = $1/648$

i) Book as a noun

$$P[\text{noun} | (\text{book}, \text{det})] = P(\text{book} | \text{noun}) * P(\text{noun} | \text{det}) + P(\text{book} | \text{det}) * \text{prev prob}$$

$$= \frac{1}{6} * \frac{1}{6} * \frac{1}{648} = \frac{1}{3888}$$

ii) Book as a verb

$$P[\text{verb} | (\text{book}, \text{det})] = P(\text{book} | \text{verb}) * P(\text{verb} | \text{det}) + P(\text{book} | \text{det}) * \text{prev prob}$$

$$= \frac{1}{4} * 0 * \frac{1}{648} = 0$$

perform smoothing

$$= \frac{1}{4} * \frac{1}{100} * \frac{1}{648} = \frac{1}{259200}$$

$\frac{1}{3888} > \frac{1}{259200}$
∴ book is a noun

∴ result of POS tagging after applying HMM

$\langle \text{S} \rangle$	The	park	is	a	book	$\langle \text{S} \rangle$
	det	noun	verb	det	noun	

Learning the parameters

Scenario 1:

• A labelled dataset is available, with the POS category of individual words in a corpus

How to do → Parameters can be directly estimated using maximum likelihood estimate from the labelled dataset.

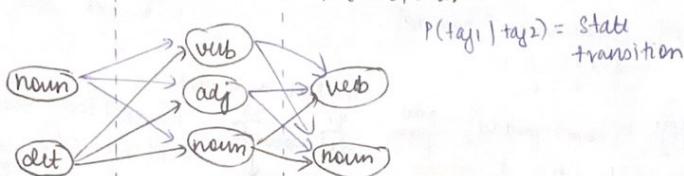
Scenario 2:

• Only the corpus is available, but not labelled with the POS categories

How to do → Baum-Welch algorithm is used to estimate the parameters of the hidden markov model.

"The light book"

$$\begin{aligned}
 P(\text{the} | \text{det}) &= 0.3 & P(\text{verb} | \text{det}) &= 0.0000 \\
 P(\text{the} | \text{noun}) &= 0.1 & P(\text{noun} | \text{det}) &= 0.5 \\
 P(\text{light} | \text{noun}) &= 0.003 & P(\text{adj} | \text{det}) &= 0.3 \\
 P(\text{light} | \text{adj}) &= 0.002 & P(\text{noun} | \text{adj}) &= 0.2 \\
 P(\text{light} | \text{verb}) &= 0.06 & P(\text{adj} | \text{noun}) &= 0.002 \\
 P(\text{book} | \text{noun}) &= 0.003 & P(\text{noun} | \text{adj}) &= 0.2 \\
 P(\text{book} | \text{verb}) &= 0.01 & P(\text{noun} | \text{verb}) &= 0.3 \\
 P(\text{word} | \text{tag}) &= \text{Symbol emission} & P(\text{verb} | \text{noun}) &= 0.3 \\
 && P(\text{verb} | \text{adj}) &= 0.001 \\
 && P(\text{verb} | \text{verb}) &= 0.1
 \end{aligned}$$



The light book

$$\begin{aligned} \text{total tags} &= 4 \\ \text{tag prob} &= 1/4 \end{aligned}$$

I) The First Word : "The"

$$\begin{aligned}
 &\xrightarrow{\alpha_1} \xrightarrow{\alpha_2} \text{det} \\
 &\quad = P(\text{det} | \text{The}, \text{ss}) \\
 &\quad = P(\text{the} | \text{det}) * P(\text{det} | \text{ss}) \\
 &\quad = 0.3 * \cancel{1/4} \downarrow \text{prob of det} \\
 &\quad = 0.075
 \end{aligned}$$

$$\begin{aligned}
 &\xrightarrow{\alpha_1} \xrightarrow{\alpha_2} \text{noun} \\
 &\quad = P(\text{noun} | \text{the}, \text{ss}) \\
 &\quad = P(\text{the} | \text{noun}) + P(\text{noun} | \text{ss}) \\
 &\quad = 0.1 + \cancel{1/4} \downarrow \text{prob of noun} \\
 &\quad = 0.025 \\
 &\quad \cancel{0.075} > 0.025 \\
 &\quad \text{det} > \text{noun} \\
 &\quad \therefore \text{The} = \text{det}
 \end{aligned}$$

II) The second word : "light"

$$P_{\text{prev prob}} = 0.075$$

$$\begin{aligned}
 \text{det} &\rightarrow \text{verb} \\
 &\quad = P(\text{verb} | \text{light, det}) \\
 &\quad = P(\text{light} | \text{verb}) * P(\text{verb} | \text{det}) * P_{\text{prev prob}} \\
 &\quad = 0.06 * 0.00001 * 0.075 \\
 &\quad = 0.00000045
 \end{aligned}$$

$$\begin{aligned}
 \text{det} &\rightarrow \text{adj} \\
 &\quad = P(\text{adj} | \text{light, det}) \\
 &\quad = P(\text{light} | \text{adj}) * P(\text{adj} | \text{det}) * P_{\text{prev prob}} \\
 &\quad = 0.002 * 0.3 * 0.075 \\
 &\quad = 0.000045
 \end{aligned}$$

$$\begin{aligned}
 \text{det} &\rightarrow \text{noun} \\
 &\quad = P(\text{noun} | \text{light, det}) \\
 &\quad = P(\text{light} | \text{noun}) * P(\text{noun} | \text{det}) * P_{\text{prev prob}} \\
 &\quad = 0.003 * 0.5 * 0.075 \\
 &\quad = 0.0001425
 \end{aligned}$$

$$0.00000045 > 0.000045 > 0.00000045$$

noun > adj > verb

$\therefore \text{light} = \frac{\text{noun}}{\text{adj}}$

III) The third word : "book"

$$P_{\text{prev prob}} = 0.0001125$$

$$\begin{aligned}
 &\cancel{\text{verb}} \rightarrow \text{verb} \\
 &\quad = P(\text{verb} | \text{book, noun}) \\
 &\quad = P(\text{book} | \text{verb}) * P(\text{verb} | \text{noun}) \\
 &\quad \quad \quad * P_{\text{prev prob}} \\
 &\quad = 0.01 * 0.3 * 0.0001125 \\
 &\quad = 0.00000034
 \end{aligned}$$

$$\begin{aligned}
 &\text{noun} \rightarrow \text{noun} \\
 &\quad = P(\text{noun} | \text{book, noun}) \\
 &\quad = P(\text{book} | \text{noun}) * P(\text{noun} | \text{noun}) \\
 &\quad \quad \quad + P_{\text{prev prob}} \\
 &\quad = 0.003 * 0.2 * 0.0001125 \\
 &\quad = 0.0000007
 \end{aligned}$$

$$0.0000007 > 0.00000034$$

The light book
 ↓ ↓ ↓
 det adj verb