



## Natural Language Processing Exam Previous Year Question Solve (Fattah Vai) - KUET CSE

Natural Language Processing (Khulna University of Engineering and Technology)



Scan to open on Studocu

## SECTION B

(Answer ANY THREE questions from this section in Script B)

5. a) Define Natural Language Processing (NLP). What are the major areas of research and development of NLP? (10)  
 b) What does  $n$ -gram mean? Drive the equation of calculating the probability for  $n$ -grams model. (10)  
 c) Consider the following corpus. (8)

$\langle s \rangle$  I am Sam.  $\langle /s \rangle$   
 $\langle s \rangle$  Sam I am  $\langle /s \rangle$   
 $\langle s \rangle$  I am Sam  $\langle /s \rangle$

$\langle s \rangle$  I do not like green eggs and Sam  $\langle /s \rangle$

Using a Bigram Language model with add-one smoothing, what is  $P(\text{Sam} | \text{am})$ ? Include  $\langle s \rangle$  and  $\langle /s \rangle$  in your counts just like any other token.

- d) What is absolute discounting? What is its advantages? (07)  
 6. a) What is closed class and open class of Part-of-Speech (POS)? Explain with example. (08)  
 b) Discuss about Rule-Based POS tagging. Write the ADVERBIAL-THE RULE. (12)  
 c) For Hidden Markov Model (HMM) POS Tagging, using the following formula, find the equation of calculating tag transition probabilities. (08)

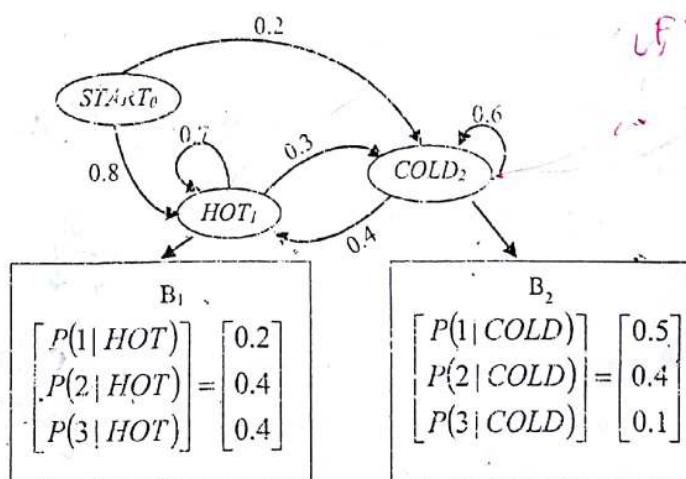
$$\hat{t}_i^n = \arg \max_{t_i^n} P(t_i^n | w_i^n)$$

- d) Consider the sentence: "Secretariat/NNP is/BEZ expected/VBN to/TO race/? Tomorrow/NR". (07)  
 The word "race" is often used as VB or NN. Given the probabilities below, find the right POS tag for the word "race".

$$P(\text{NN} | \text{TO}) = 0.00047, P(\text{VB} | \text{TO}) = 0.83, P(\text{race} | \text{NN}) = 0.00057, P(\text{race} | \text{VB}) = 0.00012,$$

$$P(\text{NN} | \text{VB}) = 0.0027, P(\text{NR} | \text{NN}) = 0.0012.$$

7. a) HMM characterized by three fundamental problems. Name and discuss about the problems. (09)  
 b) Given a sequence of ice-cream observations 313 and an HMM  $\lambda = (A, B)$  in the following figure, find the best hidden weather sequence  $Q$  (like H H H). (12)



- c) Define the term odds for logistic regression. Show that the observation should be labeled true (09) if  $\sum_{i=0}^N w_i f_i > 0$ .  
 d) Write the three-steps of Forward Algorithm. (05)  
 8. a) Name and discuss about the types of TTS. (06)  
 b) Speech Synthesis perform text to waveform mapping in two-steps. Name and discuss about the steps. Using Hourglass Metaphor. (12)  
 c) What is Homograph disambiguation? What are the problems of CMU? How does UNISYN (10)-overcome the problems of CMU?  
 d) Define text normalization. Why does text normalization important for Speech Synthesis? (07)

## N grams

1. Why would you want to predict upcoming words, or assign probabilities to sentences?
2. Define LM, N-gram, bigram, MLE, normalize, log probabilities, extrinsic and intrinsic evaluation, training set, test set, held out and development set,  
...
3. Drive the N gram probability calculation equation with example of Sam.
4. Given some sentences and count of words and size of V, draw the figures (4.1 and 4.2 in pdf page 40) and example of “I want English food”.
5. What does it mean to “fit the test set”?
6. Define Perplexity, drive equation, example
7. Define zeros situation and effect of it.
8. Define closed vocabulary and OOV, Describe the ways to train the probabilities of unknown word model.
9. Define smoothing, drive the equation for laplace smoothing with example and equation of add-k smoothing.

10. Define backoff, interpolation with equation, katz backoff with eqation, absolute discounting and equation, stupid backoff and equation.
11. Define entropy and equation, example of horse, drive the equation of Perplexity's Relation to Entropy (pdf page 56 and 57)
12. Exercises of chapter: 4.1, 4.2, 4.4, 4.7

Important for CT only: 1, 2, 3, 4, 9, 12

1: Reason for predicting upcoming words, or antigen probabilities to sentences:- Probabilities are essential in any of the following tasks

- To identify words in noisy ambiguous input. eg. Handwriting recognition, speech recognition.
- To perform auto correction, spell correction. eg their ~~are~~ is a blackboard. Here there is mistyped as their.
- for machine translation, eg. translating a bengali sentence to english with appropriate words. eg.
  - জানি আমি খুব খুব
  - I had eaten rice
- for augmented communication -  
eg. eye movement to an instruction,

[ Identify words  
[ Auto correction  
[ Machine translation  
[ Augmented Communication

CT-2  
NLP

## N grams

1 2 3 4 5 6 7 8 9  
10 11 12

2/ LM:- Models that assign probabilities to sequences of words are called language models or LMs.

MLE:- An intuitive way to estimate probabilities is called maximum likelihood estimation or MLE.

We get the MLE estimate for the parameters of an n-gram model by getting counts from a corpus, and normalizing the counts so that they lie between 0 and 1.

Normalizing:- For probabilistic models, normalizing means dividing by some total count so that the resulting probabilities fall legally between 0 and 1.

Log probability:- We always represent and compute language model probabilities in log format that is known as log probabilities.

Extrinsic Evaluation:- The best way to evaluate the performance of a language model is to embed it in an application and measure how much the application improves. Such end-to-end evaluation is called extrinsic evaluation.

Intrinsic Evaluation:- An intrinsic evaluation metric is one that measures the quality of a model independent of any application.

Training Set:- The probabilities of an n-gram model come from the corpus it is trained on known as training set or training corpus.

Test Set:- The quality of an n-gram model is measured by its performance on some unseen data, which is called test set or test corpus.

Held out set:- Test sets and other datasets that are not in our training sets are called held out corpora because we hold them out from the training data.

Development set:- A fresh test set that is truly unseen <sup>sometimes</sup> is needed and in such cases, we call the initial test set the development test set or devset.

### 3. N-Gram probability calculation equation:-

Let "Sam is a good" is an incomplete sentence. We want to know the probability of occurring boy after this sentence part

$$P(\text{boy} | \text{Sam is a good}) = \frac{c(\text{Sam is a good boy})}{c(\text{Sam is a good})}$$

Now for the full sequence of word "Sam is a good boy" we can use the chain rule -

$$P(x_1, x_2, \dots, x_n) = P(x_1) P(x_2|x_1) P(x_3|x_1, x_2) \dots P(x_n|x_1^{n-1})$$

for words ( $w_1, w_2, w_3, \dots, w_n$ )

$$P(w_1, w_2, \dots, w_n) = P(w_1) P(w_2|w_1) P(w_3|w_1^2) \dots P(w_n|w_1^{n-1})$$

$$= \prod_{i=1}^n P(w_i|w_1^{i-1})$$

$$\Rightarrow P(w_1^n) = \prod_{i=1}^n P(w_i|w_1^{i-1}) \quad \text{--- ①}$$

But calculating probability sequences are tiresome, rather we can consider only previous one or two words. e.g.  $P(\text{boy} | \text{good})$

That implies,

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-1}) \text{ for bigram}$$

This is called markov assumption

For N-gram model,

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$$

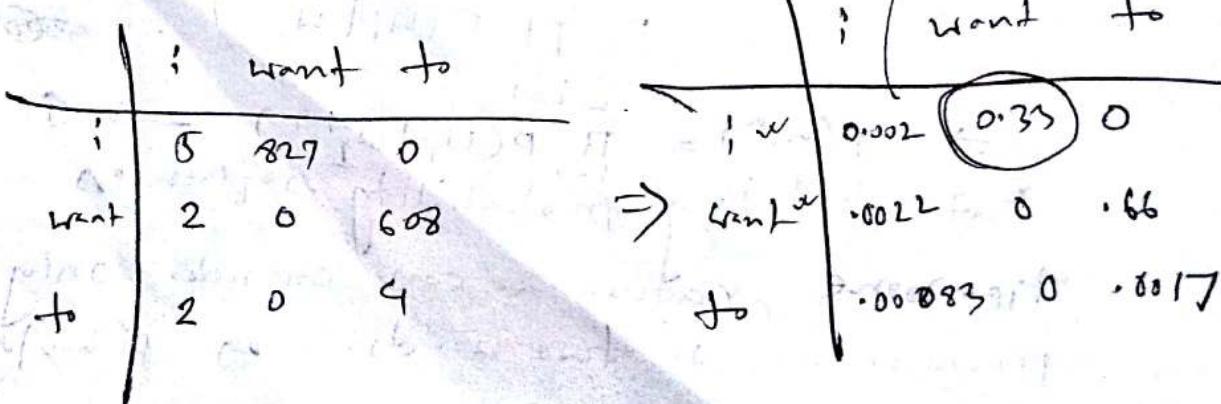
Substituting ② in ①,

$$P(w_i^n) \approx \prod_{i=1}^n P(w_i | w_{i-1})$$

Q. I want to eat Chinese food (lunch spend)

2533 927 2417 746 158 1093 341 278

$$V = 1446$$



$$p(i1 \leftrightarrow) = 0.25$$

$$p(\text{IS} | \text{foo2}) =$$

smoothing

$$P_{\text{Laplace}}^*(u_n | u_{n-1}) = \frac{c(u_{n-1}, u_n) + 1}{c(u_{n-1}) + v}$$

$$\begin{aligned} & \left. \begin{array}{l} : \text{want to} \\ : u_{n-1} \\ + \end{array} \right\} \Rightarrow \begin{array}{l} : \text{want to} \\ \frac{5+1}{2533+1446} \\ : u_{n-1} \\ + \end{array} \right\} \end{aligned}$$

5 What does it mean to "fit the test set"?

Ans:- "fit the test set" means :- whichever model assigns a higher probability to the test set - meaning if more accurately predicts the test set - is a better model. Given two probabilistic models, the better model is the one that has a tighter fit to the test data or that better predicts the details of the test data, and hence will assign a higher probability to the test data.

6 Define perplexity, derive equation, example.

Ans:- Perplexity :- ~~The~~ The perplexity (PP) of a language model on a test set is the inverse probability of the test set, normalized by the number of words.

for a test set  $w = w_1, w_2, \dots, w_N$ :

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

We can use the chain rule to expand the probability of  $w$ :

$$\text{PP}(w) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

Thus, if we are computing the perplexity of  $w$  with a bigram language model, we get :

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

The higher the 'gram' the ~~better~~ lesser the perplexity ~~which is a better approach~~

for example, considering the task of recognizing the digits (zero, one, ... nine) in english each having probability  $P = \frac{1}{10}$ , imagine a string of digits of length  $N$ .

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

$$= \left(\left(\frac{1}{10}\right)\right)^{-\frac{1}{N}}$$

$$= 10$$

$$(w_1, w_2, \dots, w_N) \rightarrow 10^{-N}$$

$$\text{if } N = 1000$$

$$\text{So } 10^{-1000} = 10^{-1000}$$

7. Define zeros situation and effects of it.

Ans:- Zeros are things that don't even occur in the training set but do occur in the test set.

Effects of zeros situation are the following:-

- First, their presence means we are under-estimating the probability of all sorts of words that might occur, which will hurt the performance of any application we want to run on this data.
- Second, if the probability of any word in the test set is zero, the entire probability of the test set is 0. By definition, perplexity is based on the inverse probability of the test set. Thus if some words have zero probability, we can't compute perplexity at all, since we can't divide by 0!

8. Define closed vocabulary and OOV. Describe the ways to train the probabilities of unknown word model.

Ans:- Closed vocabulary :- A closed vocabulary system is where the test set can only contain words from a lexicon, and there will be no unknown words.

OOV :- In some cases we have to deal with words we haven't seen before, which we'll call unknown words or out of vocabulary (OOV) words.

There are two common ways to train the probabilities of the unknown word model <UNK>

• The first one is to turn the problem back into a closed vocabulary one by choosing a fixed vocabulary in advance:

1. Choose a vocabulary that is fixed in advance.
2. Convert the training set any word that is not in this set to the unknown word token <sup><UNK></sup> in a text normalization step.
3. Estimate the probabilities for <UNK> from ~~its~~ counts just like any other regular word in the training set.

- The second alternative, in situations where we don't have a prior vocabulary in advance is to create such a vocabulary implicitly, replacing words in ~~the~~ the training data by `<UNK>` based on their frequency.

9. Smoothing :- To keep a language model from assigning zero probability to unseen events, we'll have to shave off a bit of probability mass from some more frequent events and give it to the events we've never seen. This modification is called smoothing or discounting.

Derivation of Laplace Smoothing equation:-

$$P(w_i) = \frac{c_i}{N}$$

Here,  $N =$  No. of tokens

$c_i$  = frequency of ~~the~~ <sup>word</sup>.

Add one to each count for laplace smoothing. For adjustment, we consider extra  $V$  observations in denominator.

$$P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V}$$

To make it easier, we define an adjusted count ~~c\*~~  $c^*$  normalizing by  $N$  where normalization factor is ~~N~~  $\frac{N}{N+V}$

$$c^* = (c_i + 1) \frac{N}{N + V}$$

Now for n-gram

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n)}{C(w_{n-1})}$$

$$P_{\text{Laplace}}^*(w_n | w_{n-1}) = \frac{C(w_{n-1} w_n) + 1}{\sum_w (C(w_{n-1}, w) + 1)}$$
$$= \frac{C(w_{n-1}, w_n) + 1}{C(w_{n-1}) + V}$$

Count can be replaced by  ~~$C(w_{n-1}, w_n)$~~

$$C^*(w_{n-1}, w_n) = \frac{[C(w_{n-1}, w_n) + 1] \times C(w_{n-1})}{C(w_{n-1}) + V}$$

Example for ~~smooth~~ (apply smoothing count)

	I	want	to	eat
I	5	827	0	9
want	2	0	608	1
to	2	0	4	686
eat	0	0	2	0

	I	want	to	eat
I	6	828	1	10
want	3	1	609	2
to	3	1	5	687
eat	1	1	3	1

Add k smoothing

$$P_{\text{Add-K}}^*(w_n | u_{n-1}) = \frac{c(w_{n-1}, w_n) + k}{c(w_{n-1}) + kV}$$

Comparing priors with observed probability of -  
• Stationary distribution of observations  
• Non-Homogeneous transition probabilities

• uses HMMs with scores also varying at each time step  
• can't be calculated directly due to initial condition  
• we resort to an auxiliary variable  
• instead of  $P(w_n | u_{n-1})$  calculate  $P(w_n | u_{n-1}, p)$   
• where  $p$  is a parameter of the model  
•  $P(w_n | u_{n-1}, p) = \sum_{w_{n-1}} P(w_n | w_{n-1}, p) P(w_{n-1} | u_{n-1}, p)$   
•  $P(w_n | w_{n-1}, p) = \frac{c(w_{n-1}, w_n) + k}{c(w_{n-1}) + kV}$   
•  $P(w_{n-1} | u_{n-1}, p) = \frac{\sum_{w_n} P(w_n | w_{n-1}, p)}{\sum_{w_{n-1}} \sum_{w_n} P(w_n | w_{n-1}, p)}$

Exercise

3.1

General Case.

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{c(w_{n-N+1}^{n-1} | w_n)}{c(w_{n-N+1}^{n-1})}$$

for trigrams,

$$P(am | <s>I) = \frac{c(<s>I am)}{c(<s>I)} = \frac{1}{2}$$

$$P(Sam | I am) = \frac{1}{2}$$

$$P(I | <s>Sam) = \frac{1}{2}$$

$$P(I | <s>Sam) = \frac{1}{1} = 1$$

3.2

$P(I \text{ want chinese food})$

$$= P(want | I) P(chinese | want) P(food | chinese)$$

3.4

$$P(\text{sam|am}) = \frac{c(\text{am sam})}{c(\text{am})}$$

$$c(\text{am sam}) = 2$$

$$c(\text{am}) = 3$$

$$v = 11$$

Applying  $\lambda=1$  smoothing

$$P(\text{sam|am}) = \frac{c(\text{am sam})+1}{c(\text{am})+v}$$

$$\approx \frac{3}{14} = 0.21428$$

3.7

$$P(\text{Sam|am}) = \lambda_1 P(\text{sam|am}) + \lambda_2 P(\text{sam})$$

$$= \lambda_1 \times \frac{c(\text{am sam})}{c(\text{am})} + \lambda_2 \times \frac{c(\text{Sam})}{N}$$

$$= \frac{1}{2} \times \frac{2}{3} + \frac{1}{2} \times \frac{9}{25}$$

$$= 0.4133$$

④ Absolute discounting: when we shave off a portion  $\frac{3+11}{14}$  a portion

(called d) from non-zero occurrences and give

that to zero or lower occurrence, is called ~

Advantages  $P(w_i|w_{i-1}) = \frac{c(w_{i-1}, w_i)}{\sum c(w_{i-1}, v)} + \gamma(w_{i-1}) P(w_i)$

- Increases the smaller counts
- Makes more accurate result
- More reliable result.

6@ Closed class: A class where it has a fixed membership. For example: Preposition, Conjunction they have fixed numbers of members. There is very little chance to add new words in them.

Open class: A class where there is a chance of inserting new members. For example: In English, there are 4 classes where new words can be included.

- i) Noun: Any naming thing ...
- ii) Adjective: Talks about quality
- iii) Verb: Action - - -
- iv) Adverb: modifies ... =

## POS tagging

2018  
6.(a)

- Closed class of Part of Speech: Closed class of POS are those having relatively fixed membership.

Example -

- Prepositions , because new prepositions are rarely coined.
- function words : grammatical words like is, and or you, which tend to be very short, occur frequently, and play an important role in grammar .

- Open class of Part of Speech: Open class of POS are those not having fixed membership and continually coined or borrowed from other languages.

Example -

- four major open classes occurring in the languages of the world : nouns, verbs , adjectives and adverbs .

Q What is POS tagging? Write short notes on its significance.

Ans. POS tagging:- POS tagging is the process of marking up a word in a corpus to a corresponding part of a speech tag, based on its context and definition.

Importance of POS tagging:-

- The significance of POS for language processing is that it gives a significant amount of information about the words and its neighbours.
- POS can be used in stemming for Information Retrieval (IR), since
  - Knowing a word's POS can help tell us which morphological affixes it can take
  - They can help an IR application by helping select out nouns or other important words from a document.

Importance of POS tagging:-

- Speech recognition
- NL Parsing
- Information Retrieval.

2018  
6(6)

## Rule-based POS tagging:-

- Involve a large database of hand-written disambiguation rule specifying, for example, that an ambiguous word is a noun rather than a verb if it follows a determiner.
- In Rule based pos tagging the earliest algorithms for automatically assigning POS were based on a two stage architecture
  - first, use a dictionary to assign each word a list of potential POS.
  - Second, use large lists of hand written disambiguation rules to narrow down this list to a single POS for each word.

## The ADVERBIAL-THT Rule:-

Given input: "that"

if  
 $(+1 \text{ A/ADV/QUANT})$ ; /\* if next word is adj, adverb, or quantity \*/  
 $(+2 \text{ SENT-LIM})$ ; /\* and following which is a sentence boundary \*/  
 $(\text{NOT-1 SVOC/A})$ ; /\* consider 'which' allows adj as object complement \*/

then eliminate non-~~ADV~~ ADV tags

else eliminate ADV tags

Example:-

I can't wait that long. - "that" used as adverb.

I consider that odd. - "that" used as determiner.

\* Stochastic Taggers:-

Resolve tagging ambiguities by using a training corpus to count the probability of a given word having a given tag in a given context.

## Hidden Markov Model

Markov Chain: A weighted automaton in which the input sequence uniquely determines which states the automaton will go through.

A markov chain is specified by the following components :

- $Q = q_1 q_2 \dots q_N$  a set of  $N$  states
- $A = a_{01} a_{02} \dots a_{ni} \dots a_{nm}$  a transition probability matrix  $A$ , each  $a_{ij}$  representing the probability of moving from state  $i$  to state  $j$ . Note that
$$\sum_{j=1}^n a_{ij} = 1 \quad \forall i$$
- $q_0, q_f$  a special start state and end state which are not associated with observations.

# ~~HMM POS Tagging~~ NLP CT #3 (fall 2021)

1. for HMM POS Tagging, using the following formula, find the equation of calculating tag transition probabilities. ③

$$\hat{t}_i^n = \arg \max_{t_i^n} P(t_i^n | w_i^n)$$

2. Consider the sentence : "secretariat/NNP is/VBZ expected/VBN to/TO race/VB? Tomorrow/NR" The word "race" is often used as VB or NN. Given the probabilities below, find the right POS tag for the word "race". ⑦

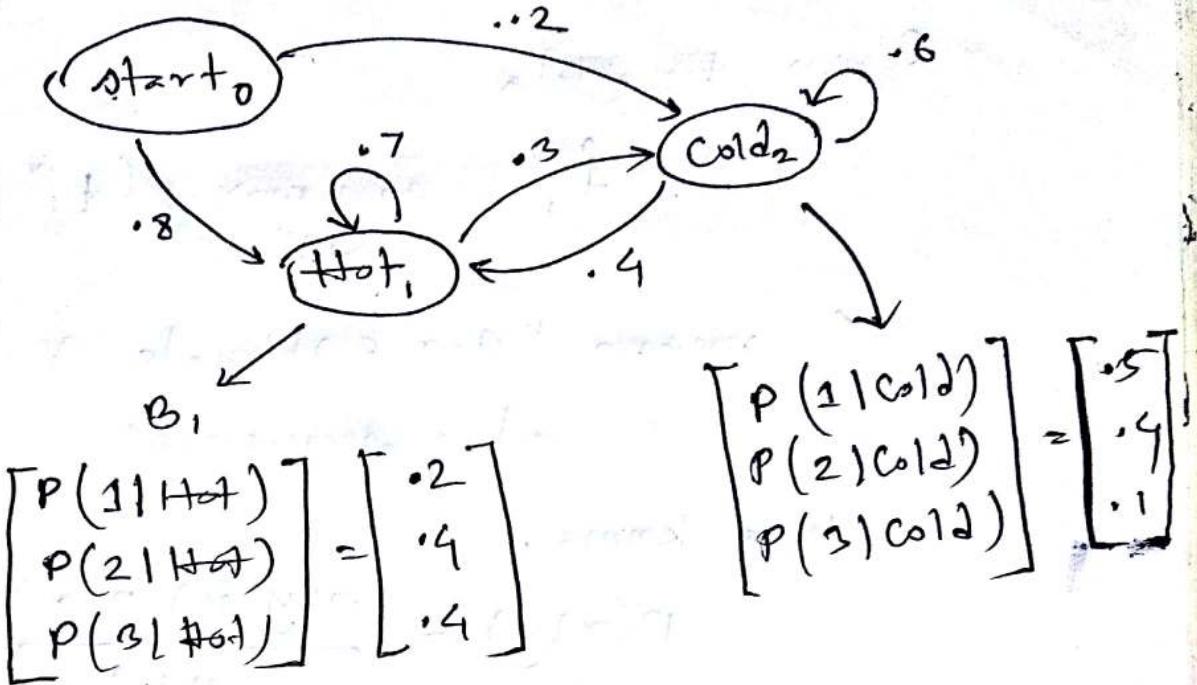
$$P(NN|T_0) = 0.00047, P(VB|T_0) = 0.83$$

$$P(race|NN) = 0.00057, P(race|VB) = 0.00012$$

$$P(NN|VB) = 0.0027, P(NR|NN) = 0.0012$$

- ✓ 3. HMM characterized by three fundamental problems. Name and Discuss about the problems. ⑨

- ✓ 4. Given a sequence of ice-cream observations 313 and an HMM  $\lambda = (A, B)$  in the following figure. find the best hidden weather sequence  $Q$  (like H-H-H) ⑫  
↓  
viterbi



✓ 5. Define the term odds for logistic regression.  
Show that the observation should be labeled true if  $\sum_{i=0}^N w_i f_i > 0$ . ⑤

6. Write the three steps of forward algorithm. ⑤

I-Am:

Given Formula

$$\hat{t}_i^n = \arg \max_{t_i^n} p(t_i^n | w_i^n)$$

$\hat{\cdot}$  mean "our estimate of the correct tag sequence"

We know,

$$p(x_i|y) = \frac{p(y|x_i) p(x_i)}{p(y)}$$

$$\therefore \hat{t}_i^n = \arg \max_{t_i^n} \frac{p(w_i^n | t_i^n) p(t_i^n)}{p(w_i^n)}$$

But  $p(w_i^n)$  doesn't change for each tag sequence

$$\therefore \hat{t}_i^n = \arg \max_{t_i^n} \underbrace{p(w_i^n | t_i^n)}_{\text{likelihood}} \underbrace{p(t_i^n)}_{\text{prior probability}}$$

Assumption 1: The probability of a word appearing is dependent only on its own part of speech tag.

$$p(w_i^n | t_i^n) \approx \prod_{i=1}^n p(w_i^n | t_i^n)$$

Assumption 2: The probability of a tag appearing is dependent only on the previous tag.

$$P(t_i^n) \propto \prod_{i=1}^n P(t_i | t_{i-1})$$

$$\hat{t}_i^n = \arg \max_{t_i^n} P(t_i^n | w_i^n) \approx \arg \max_{t_i^n} P(w_i | t_i) P(t_i | t_{i-1})$$

The tag transition probabilities:

$$P(t_i | t_{i-1}) = \frac{c(t_{i-1}, t_i)}{c(t_{i-1})}$$

The word likelihood probabilities:

$$P(w_i | t_i) = \frac{c(t_i, w_i)}{c(t_i)}$$

2. Ans.

To - verb noun

To - noun noun

for verb,  $P(VB|TO) P(NR|VB) P(race|VB) = .0000027$

for noun,  $P(NN|TO) P(NR|NN) P(race|NN) = .0000000032$

$$.00000027 > .0000000032$$

∴ The right POS tag for the word "race" is "VB".

5. Ans:

Odds:- In logistic regression, odds are defined as the ratio of the probability of success and the probability of failure.

$$\frac{p(y = \text{true} | x)}{1 - p(y = \text{true} | x)} = w \cdot f \quad (\text{dot product of weight vector and feature vector})$$

According to logit function.

$$\ln \left( \frac{p(y = \text{true} | x)}{1 - p(y = \text{true} | x)} \right) = w \cdot f$$

$$\Rightarrow \frac{p(y = \text{true} | x)}{1 - p(y = \text{true} | x)} = e^{w \cdot f}$$

$$\Rightarrow p(y = \text{true} | x) = (1 - p(y = \text{true} | x)) e^{w \cdot f}$$

$$\Rightarrow p(y = \text{true} | x) = e^{w \cdot f} - p(y = \text{true} | x) e^{w \cdot f}$$

$$\Rightarrow p(y = \text{true} | x) (1 + e^{w \cdot f}) = e^{w \cdot f}$$

$$\therefore p(y = \text{true} | x) = \frac{e^{w \cdot f}}{1 + e^{w \cdot f}}$$

~~+~~ ~~\*~~ ~~+~~ ~~\*~~

$$\begin{aligned} p(y \neq \text{false} | x) &= 1 - \frac{e^{w_f}}{1 + e^{w_f}} \\ &= \frac{1 + e^{w_f} - e^{w_f}}{1 + e^{w_f}} \\ &= \frac{1}{1 + e^{-w_f}} \end{aligned}$$

Our observation ~~will~~ should be labeled  
'true' is

$$p(y = \text{true} | x) > p(y = \text{false} | x)$$

$$\Rightarrow \frac{p(y = \text{true} | x)}{p(y = \text{false} | x)} > 1$$

$$\Rightarrow \frac{e^{w_f}}{1 + e^{w_f}} > 1$$

$$\Rightarrow e^{w_f} > 1$$

$$\Rightarrow w_f > 0$$

$$\Rightarrow \sum_{i=0}^n w_i b_i > 0$$

[Proved]

3. Ans.

HMM characterized by three fundamental problems, that are

- Problem 1 - Computing Likelihood
- Problem 2 - Decoding
- Problem 3 - Learning

Computing Likelihood :- Given an HMM  $\lambda = (A, B)$  and an observation sequence  $O$ , determine the likelihood  $\phi(O|\lambda)$ .

The most common algorithm for computing Likelihood in HMM is "The Forward Algorithm".

Decoding :- Given an observation sequence  $O$  and an HMM  $\lambda = (A, B)$ , discover the best hidden state sequence  $Q$ .

The most common decoding algorithm for HMMs is "The viterbi Algorithm".

Learning :- Given an observation sequence  $O$  and the set of states in the HMM, learn the HMM parameters  $A$  and  $B$ .

The most common learning algorithm for HMM is "The forward-Backward Algorithm".

6. Ans:

Three steps of forward Algorithm are -

1. Initialization

$$\alpha_1(j) = a_0, b_j(0_1) \quad 1 \leq j \leq N$$

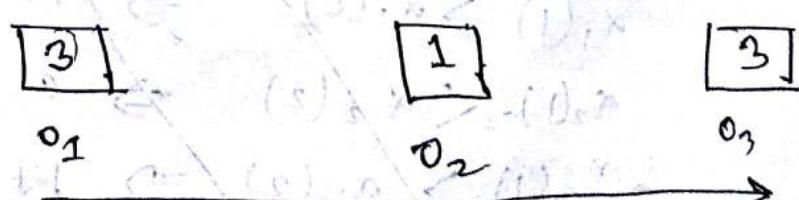
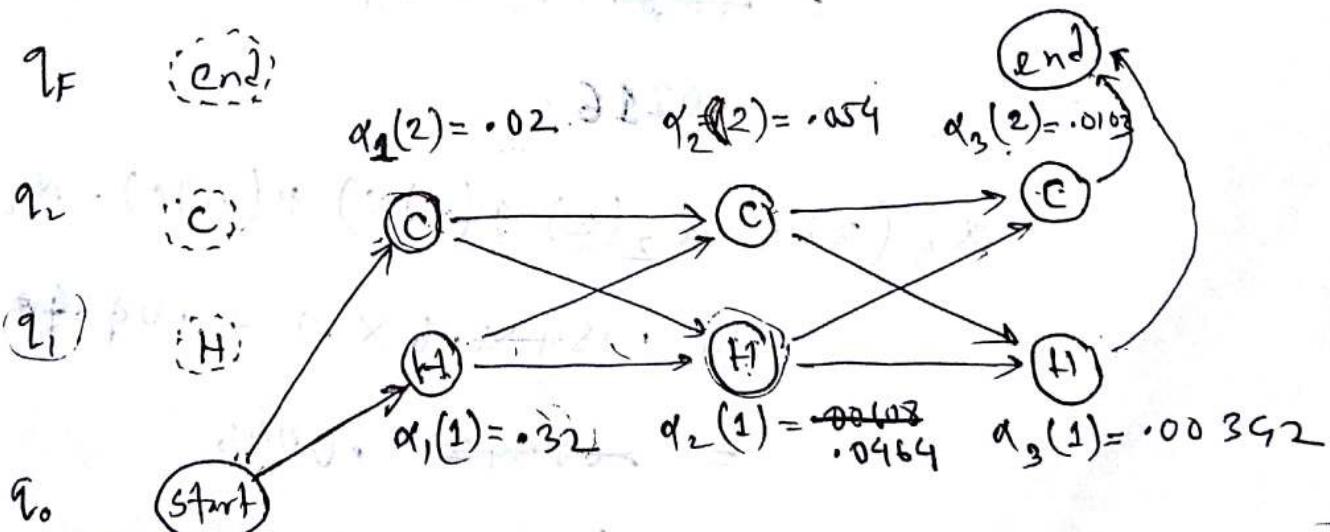
2. Recursion (Since state 0 and F are not emitting)

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(0_t); \quad 1 \leq j \leq N, 1 < t \leq T$$

3. Termination

$$\Phi(O|z) = \alpha_T(z_F) = \sum_{i=1}^N \alpha_T(i) a_{iF}$$

4. Ans: Using forward Algorithm:-



$$\alpha_1(1) = P(H|start) * P(3|H) = .8 \times .4 = .32$$

$$\alpha_1(2) = P(c|start) * P(3|C) = .2 \times .1 = .02$$

$$\begin{aligned}\alpha_2(1) &= \alpha_1(1) * P(H|H) * P(1|H) + \alpha_1(2) * P(H|C) * P(1|H) \\ &= .32 \times .7 \times .2 + .02 \times .4 \times .2\end{aligned}$$

$$\begin{aligned}\alpha_2(2) &= \alpha_1(2) * P(C|C) * P(1|C) + \alpha_1(1) * P(C|H) * P(1|H) \\ &= .02 \times .6 \times .5 + .32 \times .3 \times .5 \\ &= .054\end{aligned}$$

$$\begin{aligned}\alpha_3(1) &= \alpha_2(1) * P(H|H) * P(3|H) + \alpha_2(2) * P(H|C) * P(3|H) \\ &= .0964 \times .7 \times .4 + .054 \times .4 \times .4\end{aligned}$$

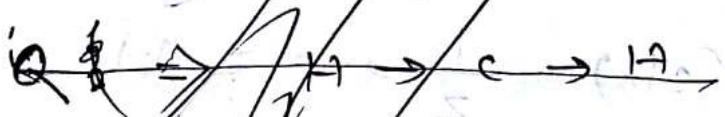
$$= .0216$$

$$\begin{aligned}\alpha_3(2) &= \alpha_2(2) * P(C|C) * P(3|C) + \alpha_2(1) * P(C|H) * P(3|C) \\ &= .054 \times .6 \times .1 + .0964 \times .3 \times .1 \\ &= .0046\end{aligned}$$

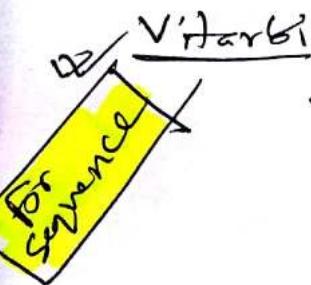
$$\begin{aligned}\alpha_1(1) &> \alpha_1(2) \Rightarrow H \\ \alpha_2(1) &< \alpha_2(2) \Rightarrow C \\ \alpha_3(1) &> \alpha_3(2) \Rightarrow H\end{aligned}$$

proper parse tree.  
question answering? Explain.

$\therefore$  Best hidden weather sequence to



$$\text{if } \alpha_1(1) > \alpha_1(2) \quad Q = HCH \\ P(O|Q) = \alpha_1(1) \times p(\text{end}/14) + \alpha_1(2) \times p(\text{end})$$



$$v_1(1) = .8 \times .4 = .32$$

$$v_1(2) = .2 \times .1 = .02$$

$$v_2(1) = \max (.32 \times .7 \times .2, .02 \times .4 \times .2) \\ = .0448$$

$$v_2(2) = \max (.32 \times .3 \times .5, .02 \times .6 \times .5) \\ = .048$$

$$v_3(1) = \max (.0448 \times .7 \times .4, .048 \times .4 \times .4) \\ = .0125$$

$$v_3(2) = \max (.0448 \times .3 \times .1, .048 \times .6 \times .1) \\ = .003$$

$$v_1(1) > v_1(2) \Rightarrow H$$

$$v_1(1) < v_2(2) \Rightarrow C$$

$$v_3(1) > v_3(2) \Rightarrow H$$

$\therefore$  Best hidden weather sequence  $Q = HCHH$ .

## \* Maximum Entropy Modelling:-

$$P(c|x) = \frac{1}{Z} \exp \left( \sum_i w_i f_i \right)$$

$Z$  = normalizing factor

$$P(c|x) = \frac{\exp \left( \sum_{i=0}^N w_i f_i \right)}{\sum_c \exp \left( \sum_{i=0}^N w_c i f_i \right)}$$

$$\sum_c \exp \left( \sum_{i=0}^N w_c i f_i \right)$$

Ex

Given, six features  $f_1, f_2, f_3, f_4, f_5, f_6$

		$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	$f_6$
VB	f	0	1	0	1	1	0
VB	w		.8		.01	.1	
NN	f	1	0	0	0	0	1
NN	w	.8					-1.3

$$P(NN|x) = \frac{e^{.8} e^{-1.3}}{e^{.8} e^{-1.3} + e^{.8} e^{.01} e^{.1}} = .2$$

$$P(VB|x) = \frac{e^{.8} e^{.01} e^{.1}}{e^{.8} e^{-1.3} + e^{.8} e^{.01} e^{.1}} = .8$$

$\therefore P(VB|x) > P(NN|x)$ , so the right pos tag for the word "race" is "VB". next are nnnb/bd/pos

## Linear Regression

- ① regression when the output is real-valued.
- ② classification when the output is one of a discrete set of classes.

#  $y = mx + c$

$$m = \frac{n \sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

$$\checkmark = \frac{n \sum x_i y_i - \cancel{n} \cancel{\sum x_i} \cancel{\sum y_i}}{n \sum x_i^2 - (\cancel{n} \cancel{\sum x_i})^2}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$c = \frac{\sum y_i - m \sum x_i}{n} = \bar{y} - m \bar{x}$$

Example  
page 46

$x$	$y$	$xy$	$x^2$
4	0	0	16
3	1000	3000	9
2	1500	3000	4
2	6000	12000	4
1	14000	14000	1
0	18000	0	0
$\sum x = 12$		$\sum y = 40500$	$\sum x^2 = 34$
		$\sum xy = 32000$	

$n = 6$

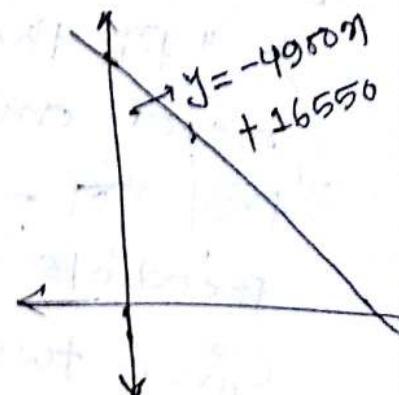
$$m = \frac{6 \times 32000 - 12 \times 40500}{6 \times 34 - (12)^2}$$

$$= -4900$$

$$c = \frac{40500 - (-4900) \times 12}{6}$$

$$= 16550$$

$$\boxed{y = -4900x + 16550}$$



## Speech Synthesis

2018  
7.(a)

### Types of TTS :-

There are three types of TTS (Text to Speech).

- Formant
- Articulatory
- Concatenative
- Articulatory Synthesis : Model movements of articulators and acoustics of vocal tract.
- Formant Synthesis : start with acoustics, create rules/ filters to create each formant.
- Concatenative synthesis : Use databases of stored speech to assemble new utterances.

A full system needs to go all the way from random text to sound.

### \* Concatenative synthesis :-

- Concatenative synthesis is a technique for synthesizing sounds by concatenating short samples of recorded sound (called units).
- The duration of the units is not strictly defined and may vary according to the implementation, roughly in the range of 10 milliseconds up to 1 second.
- All current commercial systems are concatenative synthesis.

2018  
8.(b)

Speech synthesis perform text to wave form mapping in two following steps :-

- Step 1 (Text analysis) : convert the input text into a phonemic internal representation.

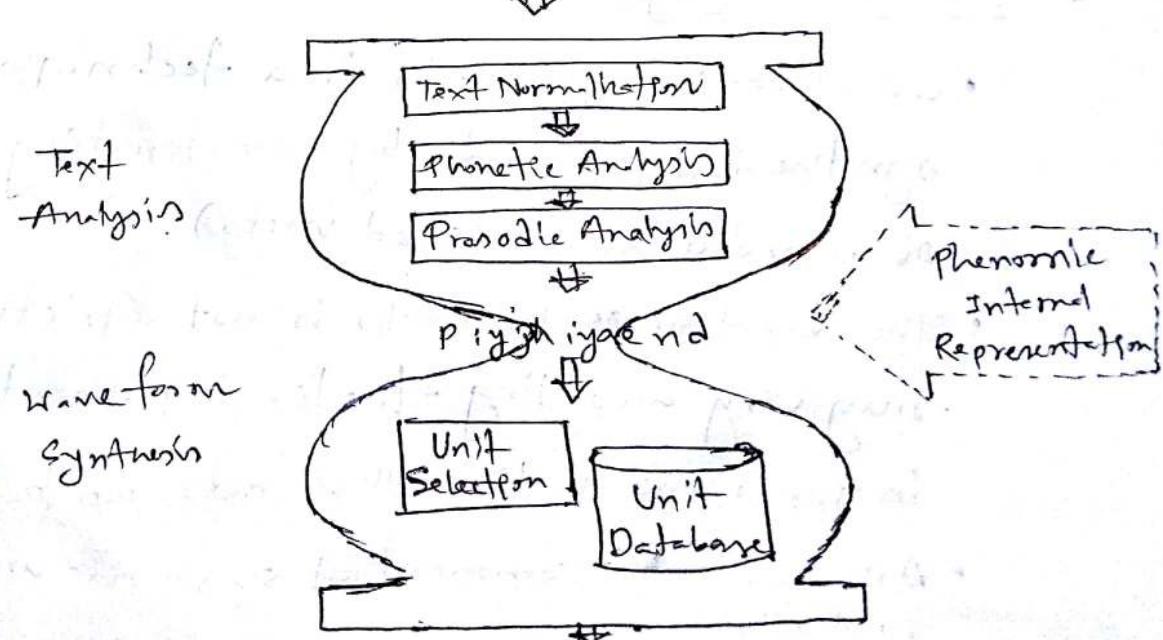
P	G	AND	E	*
p iy	jh iy	ae n d	iy	...

- Step 2 (Wave-form synthesis) : convert this internal representation into a waveform



According to the homoglass metaphor

PG&E will file schedules on April 20



2018  
8(c)

Homograph Disambiguation :- Homographs are words that are spelled the same but have different meanings.

The process of making it disambiguous is known as homograph disambiguation. e.g.

Do you live near a zoo with live animals.

/ɪn vɪ/

/læv vɪ/

It is not a huge problem, but still important to make it disambiguous.

Problems of CMU :-

- Flaw errors
- Only American pronunciation
- No syllable boundaries
- Doesn't tell us which pronunciation to use for which homophones
  - no POS tags
- Doesn't distinguish case
  - The word US has 2 pronunciations  
[Aʌs] and [y uəs EHs]

UNISYN solve problems of CMU in the following way:-

- Has syllabification, stress, some morphological boundaries.
- Pronunciation can be read off in
  - General American
  - RP British
  - Australia
  - Etc

2018  
8.(d)

Text normalization :- Text normalization is the process of transforming text into a single canonical form that it might not have had before

paragraph → sentence → tokens.

Importance of text normalization for Speech Synthesis :-

- Identify tokens in text
- Chunk tokens into reasonably sized sections
- Map tokens to words
- Identify types of words.

2013  
5.(a)

NLP :- NLP is a field of computer science, artificial intelligence, and linguistics concerned with the interactions between computers and human languages.

### Major areas of R&D in NLP :-

- Text processing
  - Processing of raw text
- Morphological Analysis - Study of word formation
  - Washing = Wash + ing
  - Running = Run + ing
- POS tagging
  - >> text = word\_tokenize("today is a beautiful day")
  - >> nltk.pos\_tag(text)

Output : [ ('today', 'NN'), ('is', 'PRP'), ('a', 'DT'), ('beautiful', 'JJ'), ('day', 'NN') ]

- Machine translation
  - Translating content in one natural language to another natural language.

Ex: આજે એ ચાંદ - I eat rice.

- Parsing
  - Identifying sense structure
- Text to speech
  - Converting electronic text to digital speech
- Automatic speech recognition
  - Automatic transcription of spoken content to electronic text.
- Speech-to-speech translation
  - Translating spoken content from one language to another in realtime or offline.

### Industrial Application of NLP:-

- Search engines
- fraud detection
- Information extraction
- Collaborative filtering
- Computational advertising
- Advanced text Editors
- Sentiment analysis
- Opinion mining

# Backoff: We use the trigram if the evidence is sufficient. Otherwise we use the bigram, otherwise the unigram. So if the higher order n-gram have zero evidence then we will "back-off" to the lower-order n-gram.

Interpolation: Interpolation is the method where we mix the probability estimates from all the n-gram estimators, weighing and combining the trigram, bigram, and unigram counts.

A simple linear interpolation

$$\hat{P}(w_n | w_{n-2} w_{n-1}) = \lambda_1 P(w_n | w_{n-2} w_{n-1}) + \lambda_2 P(w_n | w_{n-1}) + \lambda_3 P(w_n)$$

such that the  $\lambda$ 's sum to 1:

$$\sum_i \lambda_i = 1.$$

# Katz backoff: In Katz backoff we apply a discount probability  $p^*$  if we've seen this n-gram before. Otherwise, we decrease back off to the Katz probability for the shorter history (n-1)-gram, the prob. for a backoff n-gram  $P_{BO}$  is

$$P_{BO}(w_n | w_{n-N+1}^{n-1}) = \begin{cases} p^*(w_n | w_{n-N+1}^{n-1}) & \text{if } w_{n-N+1}^{n-1} \text{ is observed} \\ \alpha(w_{n-N+1}^{n-1}) P_B(w_n | w_{n-N+1}^{n-1}) & \text{otherwise,} \end{cases}$$

# absolute discounting: Absolute discounting formalizes this by subtracting a fixed discount  $d$  from each count. The intuition is that since we have good estimate already for the very high counts, a small discount  $d$  won't affect them much. For bigram,

$$P_{AD}(w_i | w_{i-1}) = \frac{c(w_{i-1} w_i) - d}{\sum_v c(w_{i-1} v)} + \lambda(w_{i-1}) \frac{P(w_i)}{P(w_{i-1})}$$

# Katz backoff: In Katz backoff we rely on a discount probability  $p^*$  if we've seen this  $n$ -gram before. Otherwise, we ~~successively~~  
back off to the Katz probability  $\lambda$  for the shorter history  $(n-1)$ -gram, the prob.  
for a backoff  $n$ -gram  $P_{BO}$  is

$$P_{BO}(w_n | w_{n-N+1}^{n-1}) = \begin{cases} p^* (w_n | w_{n-N+1}^{n-1} \text{ is } C(w_{n-N})) & \text{if } C(w_{n-N}) \\ (\lambda(w_{n-N+1}^{n-1})) P_{BO}(w_n | w_{n-N+1}^{n-1}) & \text{otherwise,} \end{cases}$$

# absolute discounting: Absolute discounting formalizes this by subtracting a fixed discount  $d$  from each count. The intuition is that since we have good estimates already for the very high counts, a small discount  $d$  won't affect them much.

for bigram,

$$P_{AD}(w_i | w_{i-1}) = \frac{C(w_{i-1} w_i) - d}{\sum_v C(w_{i-1} v)} + \lambda(w_{i-1}) \frac{P(w_i)}{P(w_{i-1})}$$

# stupid backoff: In this algorithm there is no discounting of the higher-order prob. if a higher-order n-gram has a zero count, we simply backoff to a lower order n-gram.

$$S(w_i | w_{i-k+1}^{i-1}) = \begin{cases} \frac{\text{Count}(w_{i-k+1}^i)}{\text{Count}(w_{i-k+1}^{i-1})} & \text{if } \text{Count}(w_{i-k+1}^i) > 0 \\ S(w_i | w_{i-k+2}^{i-1}) & \text{otherwise} \end{cases}$$

II Entropy: Entropy is a measure of information. Given a random variable  $X$  over a range over whatever we are predicting and with a particular probability function, call it  $p(x)$ , the entropy of the rand. var.  $X$  is,

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

Entropy is often used to measure the

information content of a message or the

surprise value of a message sent.

# Example of horse: suppose we want to bet on a horse race but it is too far away. So we'd like to send encoded messages using binary representation. So, the horse 1 001, horse 2 010, horse 3 011, and so one with horse 8 encoded as 000, is we spent all day to bet and each horse is coded with 3 bits. On avg. we will send 3 bit. Suppose the probability of winning of the horses

horse 1	$\frac{1}{2}$	horse 5	$\frac{1}{64}$
horse 2	$\frac{1}{4}$	(1)	6. 11
horse 3	$\frac{1}{8}$	"	7. 11
horse 4	$\frac{1}{16}$	"	8. 11

So, now we will try to find the entropy of the random variable  $X$  that changes over horses gives us a

lower bound on the number of bits and its

$$\begin{aligned} H(n) &= - \sum_{i=1}^{i=8} p(i) \log p(i) \\ &= -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{8} \log \frac{1}{8} \\ &\quad - \frac{1}{16} \log \frac{1}{16} - 4 \left( \frac{1}{64} \log \frac{1}{64} \right) \\ &= 2 \text{ bits.} \end{aligned}$$

so we can create a coding system that averages 2 bits per race,

so we can encode the most likely horse with 0, remaining horses as L0, L10, L1L0, LLL100, LLLL01, L11110, L11111

if all winning prob. is same and let  $\frac{1}{8}$

then;

$$\begin{aligned} H(n) &= - \sum_{i=1}^{i=8} \frac{1}{8} \log \frac{1}{8} \\ &= -\log \frac{1}{8} = 3 \text{ bits.} \end{aligned}$$

reln with perplexity: The perplexity of a model  $p$  on a sequence of words  $\omega$  is now formally defined as the exp of this entropy:

$$\text{perplexity}(\omega) = 2^{H(\omega)}$$

$$= p($$

$$H(\omega) = -\frac{1}{N} \log p(\omega_1, \omega_2, \dots, \omega_N)$$

$$\therefore \text{perplexity}(\omega) = 2^{H(\omega)}$$

$$= p(\omega_1, \omega_2, \dots, \omega_N)^{-1/N}$$

$$= N \sqrt[N]{\frac{1}{p(\omega_1, \omega_2, \dots, \omega_N)}}$$

$$= N \sqrt[N]{\prod_{i=1}^N \frac{1}{p(\omega_i | \omega_1, \dots, \omega_{i-1})}}$$