

# COMP 6721 Applied Artificial Intelligence (Fall 2023)

## Worksheet #9: Introduction to Natural Language Processing (NLP)

**Language Model.** In Natural Language Processing (NLP), a *bigram* language model is a simple yet effective way to understand the probability of a word sequence. It calculates the likelihood of a word  $w_n$  appearing after a given word  $w_{n-1}$ . We use *Maximum Likelihood Estimation* (MLE) to determine these probabilities from a given corpus:  $P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$ , where  $C(w_{n-1}w_n)$  is the count of the occurrence of  $w_{n-1}$  followed by  $w_n$ . So, given the following corpus of three sentences:

<s> I am Sam </s>  
 <s> Sam I am </s>  
 <s> I do not like green eggs and ham </s>

compute the following bigram probabilities:

$P(I|<s>) =$  .....  $P(\text{Sam}|<s>) =$  .....  $P(\text{am}|I) =$  .....  
 $P(</s>|\text{Sam}) =$  .....  $P(\text{Sam}|\text{am}) =$  .....  $P(\text{do}|I) =$  .....

**Sentence Probability.** Given an English language model with the following *bigram* probabilities, compute the probability for the sentence “I want to eat British food”:

$P(\text{on} \text{eat}) = .16$	$P(\text{want} I) = .32$	$P(\text{eat} \text{to}) = .26$	$P(I \text{ want to eat British food})$ $=$ ..... $=$ ..... $=$ .....
$P(\text{some} \text{eat}) = .06$	$P(\text{would} I) = .29$	$P(\text{have} \text{to}) = .14$	
$P(\text{British} \text{eat}) = .001$	$P(\text{don't} I) = .08$	$P(\text{spend} \text{to}) = .09$	
...	...	...	
$P(I <s>) = .25$	$P(\text{to} \text{want}) = .65$	$P(\text{food} \text{British}) = .6$	
$P(I'd <s>) = .06$	$P(a \text{want}) = .5$	$P(\text{restaurant} \text{British}) = .15$	
$P(</s> \text{British}) = .1$	$P(</s> \text{food}) = .25$	$P(</s> \text{restaurant}) = .35$	

**Corpus Probabilities.** Given a corpus with  $|V| = 1616$  different words and a total of  $N = 10000$  bigrams:

	<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<i>Total</i>
<i>I</i>	8	1087	0	13	0	0	0		$C(I)=3437$
<i>want</i>	3	0	786	0	6	8	6		$C(\text{want})=1215$
<i>to</i>	3	0	10	860	3	0	12		$C(\text{to})=3256$
<i>eat</i>	0	0	2	0	19	2	52		$C(\text{eat})=938$
<i>Chinese</i>	2	0	0	0	0	120	1		$C(\text{Chinese})=213$
<i>food</i>	19	0	17	0	0	0	0		$C(\text{food})=1506$
<i>lunch</i>	4	0	0	0	0	1	0		$C(\text{lunch})=459$
...									...
									...
									$N=10,000$

compute the probabilities for  $P(II) =$  .....,  $P(I|I) =$  ..... and  $P(\text{lunch}|I) =$  .....

**Smoothing.** We can avoid zero probabilities by *smoothing*, here we use add-one (or *Laplace*) smoothing:

	<i>I</i>	<i>want</i>	<i>to</i>	<i>eat</i>	<i>Chinese</i>	<i>food</i>	<i>lunch</i>	...	<i>Total</i>
<i>I</i>	8+1	1087+1	1	14	1	1	1		
<i>want</i>	3+1	1	787	1	7	9	7		$C(\text{want}) +  V  = 2831$
<i>to</i>	4	1	11	861	4	1	13		$C(\text{to}) +  V  = 4872$
<i>eat</i>	1	1	23	1	20	3	53		$C(\text{eat}) +  V  = 2554$
<i>Chinese</i>	3	1	1	1	1	121	2		$C(\text{Chinese}) +  V  = 1829$
<i>food</i>	20	1	18	1	1	1	1		$C(\text{food}) +  V  = 3122$
<i>lunch</i>	5	1	1	1	1	2	1		$C(\text{lunch}) +  V  = 2075$
...									

computing the new bigram probabilities as  $P_{\text{Add1}}(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)+1}{C(w_{n-1})+|V|}$  and  $P_{\text{Add1}}(w_{n-1}w_n) = \frac{C(w_{n-1}w_n)+1}{N+B}$ , where  $B$  is the number of “bins” we added +1 to (so here,  $|V|^2$ ). Compute the new probabilities:

$P(II) =$  .....,  $P(I|I) =$  ..... and  $P(\text{lunch}|I) =$  .....

**Part-of-Speech Tagging.** Given the following lexicon, assign a *part-of-speech* (POS) tag to each word for the sentence below:

### Lexicon:

N --> flight | trip | breeze | morning // noun  
 V --> is | prefer | like // verb  
 Adj --> direct | cheapest | first // adjective  
 Pro --> me | I | you | it // pronoun  
 PN --> Chicago | United | Los Angeles // proper noun  
 D --> the | a | this // determiner  
 Prep --> from | to | in // preposition  
 Conj --> and | or | but // conjunction

I	prefer	a	direct	flight	to	Chigaco.

**Parsing.** Now, given the following context-free grammar:

### Grammar:

S --> NP VP // I + prefer United  
 NP --> Pro | PN | D N | D Adj N // I, Chicago, the morning  
 VP --> V | V NP | V NP PP // is, prefer + United,  
 PP --> Prep NP // to Chicago, to I ??

create a *parse tree* for the sentence, “*I prefer a direct flight to Chicago.*” using the POS tags you assigned above:

**Word Sense Disambiguation.** Using the following probabilities you obtained from a training corpus ( $|V|=50$ ):

- $P(\text{the}|\text{BANK1}) = (5+.5) / (30+.5V)$
- $P(\text{world}|\text{BANK1}) = (1+.5) / 55$
- $P(\text{and}|\text{BANK1}) = (1+.5) / 55$
- $P(\text{off}|\text{BANK1}) = (0+.5) / 55$
- $P(\text{Potomac}|\text{BANK1}) = (0+.5) / 55$
- $P(\text{BANK1}) = 5/7$
- $P(\text{the}|\text{BANK2}) = (3+.5) / (12 + .5V)$
- $P(\text{world}|\text{BANK2}) = (0+.5) / 37$
- $P(\text{and}|\text{BANK2}) = (0+.5) / 37$
- $P(\text{off}|\text{BANK2}) = (1+.5) / 37$
- $P(\text{Potomac}|\text{BANK2}) = (1+.5) / 37$
- $P(\text{BANK2}) = 2/7$

Using “add 0.5” smoothing as shown above, with a context window of  $\pm 3$ , find the correct sense for *bank* in the sentence, “*I like the Potomac bank*”:

1. Score(BANK1) = .....

2. Score(BANK2) = .....

Note: Words not shown in the list above have an *unsmoothed* probability of 0. Use logs.