

## NLP Sheet 2

- Missing From Sheet 1

6

5. Mention two words that have

ed, s, ing, tion, or, op...

		lemma	Stem
Same lemmas and Same Stems	Played Playing	Play Play	Play Play
Same lemmas and different Stems	lie lied	lie lie	lie li
different lemmas and Same Stems	action aced	action ace	ac ac
different lemmas & different Stems	walk run	walk run	walk run

1. Write the equation for trigram Probability estimation

$$P(w_i | w_{1:i-1}) \approx P(w_i | w_{i-2} w_{i-1}) = \frac{C(w_{i-2} w_{i-1} w_i)}{C(w_{i-2} w_{i-1})}$$

N-gram  $\rightarrow$  N-1 words

then write all non-zero Probabilities for  
3-gram

①  
⑤  
<S><S> I am Sam <IS> <S><S> Sam I am <IS>  
<S><S> I do not like green eggs and ham <IS>

$$P(I | <S><S>) = \frac{2}{3} \leftarrow \begin{matrix} C(<S><S>I) \\ C(<S><S>) \end{matrix}$$

• Count denominator 1st. ⑥

$$P(am | <S>I) = \frac{1}{2} \leftarrow \begin{matrix} C(<S>Iam) \\ C(<S>I) \end{matrix}$$

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

$$P(<IS> | am Sam) = \frac{1}{1}$$

$$P(<S> | Sam <IS>) = \frac{1}{1}$$

$$P(<S> | <IS><S>) = \frac{2}{2}$$

$$P(Sam | <S><S>) = \frac{1}{3}$$

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

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

$$P(<IS> | Iam) = \frac{1}{2}$$

$$P(<S> | am <IS>) = \frac{1}{1}$$

②  
⑤  
...

$$P(<S> | <IS><S>) = \frac{2}{2}$$

\* 8 more Probabilities, all equal to 1

- Note that by writing the Probabilities in that way, all of them are guaranteed to be non zero (the other approach is to build a count matrix with rows being all consecutive word pairs and columns being all words)  
(unique)

2. Given is a bigram Probability matrix

$$w_{i-1} \left\{ \underbrace{\left( \begin{array}{c} \phantom{0} \\ \phantom{0} \\ \phantom{0} \end{array} \right)}_{w_i} \right. \quad \cdot \text{normal bigram and add-1 smoothed version.}$$

- Compute Probability of Sentence "I want Chinese Food"

→  $\langle S \rangle$  I want Chinese Food  $\langle S \rangle$

$$\text{Recall, } P(w_{1:n}) = \prod_{i=1}^n P(w_i | w_{1:i-1}) \quad \cdot i=1 \leftrightarrow \bar{I} \text{ after } \langle S \rangle$$

a) using normal Bigram Probabilities:

$$P(w_{1:n}) = P(I | \langle S \rangle) P(\text{want} | I) P(\text{Chinese} | \text{want}) P(\text{Food} | \text{Chinese}) P(\langle S \rangle | \text{Food})$$

$$= 0.25 \times 0.33 \times 0.0065 \times 0.52 \times 0.68$$

$$= 0.00189$$

b) using add-1 Smoothed Probabilities:

$$= 0.19 \times 0.21 \times 0.0029 \times 0.052 \times 0.4$$

$$= 0.0000024$$

c) unsmoothed is higher because smoothing moves Probability mass from events with non zero Probability



to events of zero Probability. This is evident by how it adds 1 in the numerator and  $|V|$  in the denominator; this decreases any fraction with a nonzero denominator as  $|V| > 1$ .

↳ vocab. length

3. Train a bigram model on the following corpus without adding an end of sentence token

$\langle s \rangle$  a b  
 $\langle s \rangle$  b b  
 $\langle s \rangle$  b a  
 $\langle s \rangle$  a a

\* Recall that training means to compute all the needed bigram probabilities.

$$w_{i-1} = \left( \begin{matrix} \vdots \\ w_i \end{matrix} \right) \quad \text{• For all possible } w_{i-1}, w_i \quad \text{©}$$

$w_{i-1}$  {

	a	b	Sum $w_{i-1}$ →
$\langle s \rangle$	2/4	2/4	4
a	1/2	1/2	2
b	1/2	1/2	2

• First write the count in black ( $w_{i-1}, w_i$ )

• Sum for the count  $w_{i-1}$  (red)

• divide to get Probab

⇒ Find the Sum of the Probabilities of all Possible two-word Sequences

→ They are aa, ab, ba, bb

$$P(\langle s \rangle aa) = P(a|\langle s \rangle) P(a|a) = \frac{2}{4} \times \frac{1}{2} = \frac{1}{4}$$

$$P(\langle s \rangle ab) = P(a|\langle s \rangle) P(b|a) = \frac{2}{4} \times \frac{1}{2} = \frac{1}{4}$$

$$P(\langle s \rangle ba) = P(b|\langle s \rangle) P(a|b) = \frac{2}{4} \times \frac{1}{2} = \frac{1}{4}$$

$$P(\langle s \rangle bb) = P(b|\langle s \rangle) P(b|b) = \frac{2}{4} \times \frac{1}{2} = \frac{1}{4}$$

→ Their Sum is Clearly 1

⇒ The Sum of Probabilities of all Possible 3-word Sequences

$$P \begin{pmatrix} \langle s \rangle ab a \\ \langle s \rangle ba a \\ \langle s \rangle aa a \\ \langle s \rangle bb a \\ \langle s \rangle ab b \\ \langle s \rangle ba b \\ \langle s \rangle aa b \\ \langle s \rangle bb b \end{pmatrix} = \begin{pmatrix} 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \end{pmatrix} \xrightarrow{\text{Sum to}} 1$$

↳ each is a Product of 3 Probabilities, each being = to half.

\* Conclusion: the Bigram model makes all Sequences of a given length have Probabilities that Sum to 1

→ That is, it doesn't assign a single Probability distribution across all Sentences.

i.e., when we compute the Probability of a Sentence, we don't get its true Probability across all Possible Sentences.

⇒ we want the distribution across all Sentences to be 1 for this

\* This is the role of the  $\langle /s \rangle$  token

→ once we add it the matrix becomes

	a	b	$\langle /s \rangle$
$\langle s \rangle$	2/4	2/4	0
a	1/4	1/4	2/4
b	1/4	1/4	2/4

and you can show now that Sum of Probabilities of all Possible Sentences is 1

	one word	two words	three words	
P	0.5	0.25	0.125	...

$$\sum_{n=1}^{\infty} \frac{1}{2^n} = 1$$

4. We are given the following corpus

- <S> I am Sam <15>
- <S> Sam I am <15>
- <S> I am Sam <15>
- <S> I do not like green eggs and Sam <15>

• the formula we already know  $C(w_{1-(n-1):i-1} w_i) / C(w_{1-(n-1):i-1})$  is a maximum likelihood estimate.

\* Use linear Interpolation between a bigram and unigram model to compute  $P(\text{Sam I am})$

$$\hat{P}(\text{Sam I am}) = \lambda_1 P(\text{Sam I am}) + \lambda_2 P(\text{Sam})$$

• given  $\lambda_1 = \lambda_2 = \frac{1}{2}$

$$\downarrow$$

$$\frac{2}{3}$$

$\leftarrow C(\text{am})$

$$\downarrow$$

$$\frac{4}{25}$$

$\uparrow$  all words

$$= 0.413$$



5. Given a training Set of 100 numbers  
→ 91 of which are zeros  
→ 9 other digits are 1, 2, 3, ..., 9

• What is unigram Perplexity for the test Set 0 0 0 0 0 3 0 0 0 0

⇒ From the training Set, it's true that

$$P(0) = \frac{91}{100}, P(3) = \frac{1}{100}$$

$$\begin{aligned} \cdot P(w_{1:n}) &= \prod_{i=1}^n P(w_i) = \left(\frac{91}{100}\right)^9 \cdot \frac{1}{100} \\ &= \frac{91^9}{100^{10}} \end{aligned}$$

*test Set*

$$\begin{aligned} PP(w) &= \sqrt[n]{\frac{1}{P(w_{1:n})}} = \sqrt[10]{\left(\frac{91^9}{100^{10}}\right)^{-1}} \\ &= 1.725 \end{aligned}$$

$\bar{x} = \frac{1}{x}$

6. Write the Feed Forward equations of an RNN

$$h_t = P(Wx_t + Uh_{t-1})$$

$$y_t = g(Vh_t)$$

• In this example,  $W=V=U=1$  (hence  $x_t$  and  $h_t$  are scalars),  $P(x)=x$  and  $g(x)=x$  (linear units)

The equations, hence, reduce to

$$h_t = x_t + h_{t-1}$$

$$y_t = h_t$$

→ given  $X = [2, -0.5, 1]$  Find all network values. ( $h_t, y_t$  for  $t \in \{0, 1, 2\}$ )

• assume  $h_{-1} = 0$

$$x_0 = 2 \rightarrow h_0 = 2 + 0 \rightarrow y_0 = 2$$

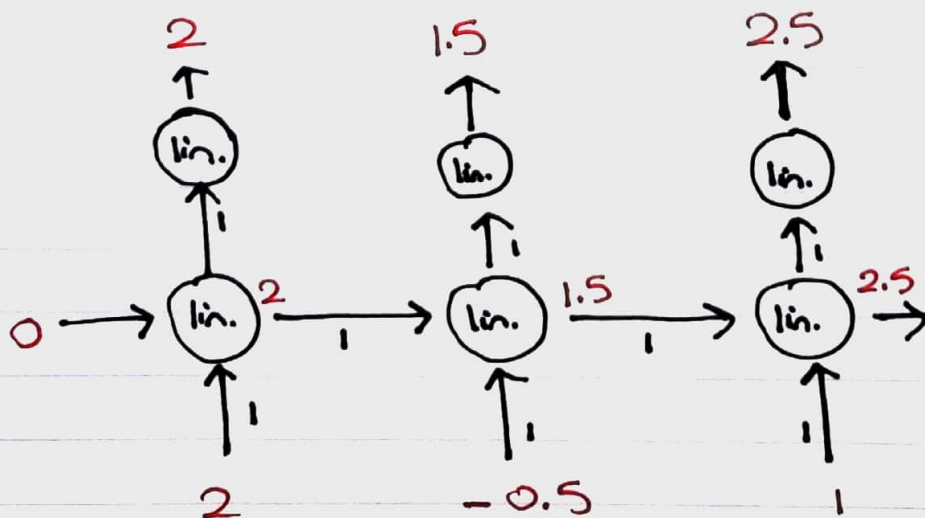
$$x_1 = -0.5 \rightarrow h_1 = -0.5 + 2 = 1.5 \rightarrow y_1 = 1.5$$

$$x_2 = 1 \rightarrow h_2 = 1 + 1.5 = 2.5 \rightarrow y_2 = 2.5$$

$$2 + -0.5 \rightarrow 1.5 + 1 \rightarrow 2.5$$

# It's learning addition (in fact, the equations by definition describe a running sum)

Unrolled Network





7.

→ given is an RNN Char. level language model. (token is a char.)

→ each token is represented by a one-hot vector  
 Vocab = {h, e, l, o}

• Use Softmax to compute final output

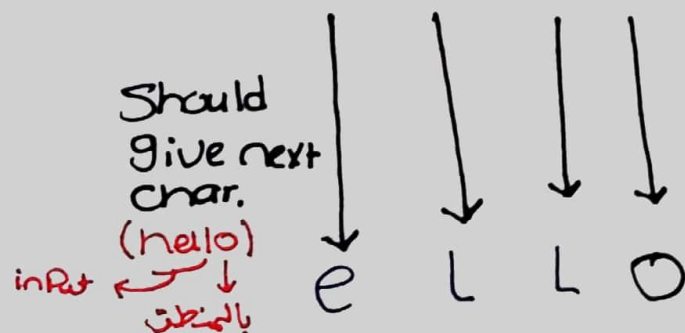
Recall,  $P(x) = \frac{e^x}{\sum_j e^{x_j}}$  ← element-wise  
 is the Softmax fun.

• Check lec. for more

→ Use it to compute the Probability dist. Vector by applying it on each of the 4 output vectors

e.g.  $P\begin{pmatrix} 1.0 \\ 2.2 \\ -3.0 \\ 4.1 \end{pmatrix} = \begin{pmatrix} e^1 \\ e^{2.2} \\ e^{-3.0} \\ e^{4.1} \end{pmatrix} \cdot \frac{1}{e^1 + e^{2.2} + e^{-3} + e^{4.1}}$

• the given input is h e l l



Recall that the output is a dist. over the vocab.

{ what the model gave

○	○	L	○
x	x	✓	✓

(Prediction is largest index in output vector)  $\begin{pmatrix} h \\ e \\ l \\ o \end{pmatrix}$

→ before or after Softmax, it's the same