## David Jose Fernandez

1

NLP HW #3

Plstore | computer) = 10 = ( value | x) | C(ui-1)

P(monitor | computer) = 4 = 2

Both words have the same bigram probability so they are both equally as likely to happen.

"refinem" mus bsnowed 2 1b) Kneser-Ney smoothing (d= 5) wicky, vi) >01

P(w; | wi-1) = max(c(vin, vi)-d,0) + 2(vi-1) = |vic(v, vi)>0|

P(store | compater) = max(3.5,0) + 1.15 ( 3+2+2+3) = .4

P(monitor 1 computer) = max (3.5,0) + 15 (2+2+2+3)=.3833

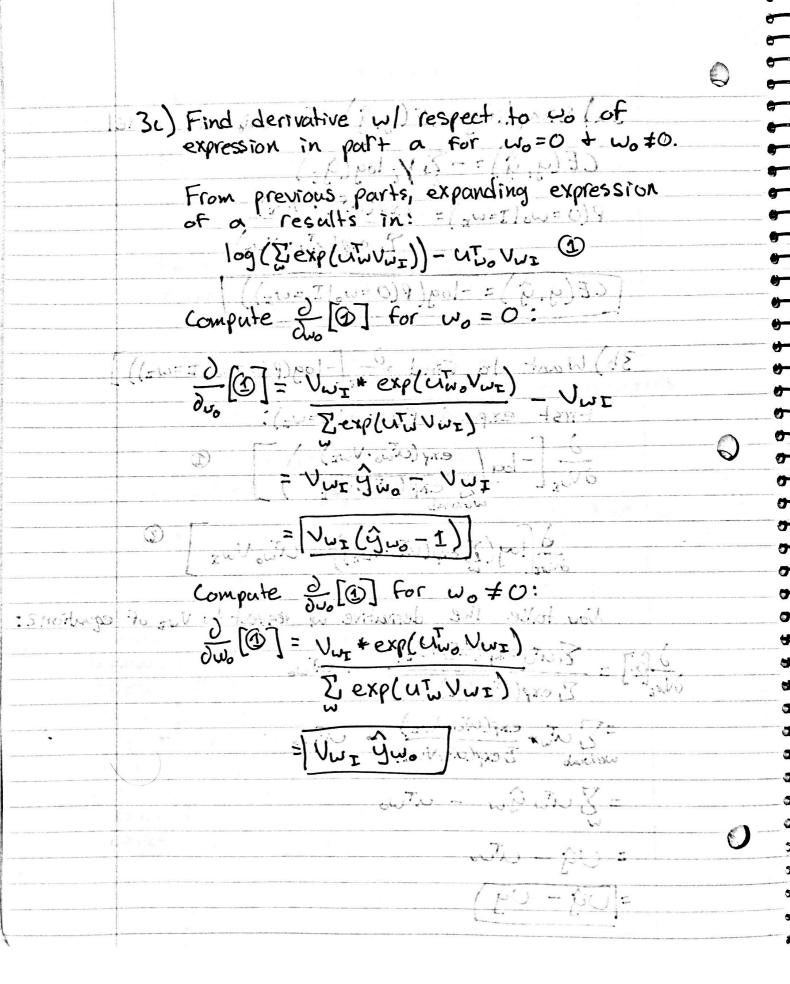
In this case, P(store | computer) > P(monitor | computer) so "store" would be schosen to complete the sentence. This makes sense since kneser-vey smoothing favors words that have higher likelihoods of being combined wil Brustions words.

GINC (32132) = 0.6786 Come(3, di) = 0, 4400 cornel 2, 3, ) = 0. 5 500

E# WH 9VU	
P(stere   compater) = $\frac{1}{10}$   $\frac{1}{10$	
= (3.4.0) (3.4.0)	
P(store   compater) =	
2/ = 1 = (19 hamo)   91012) 1	
Allomputer) = 100 - 100   rodinary)	
(Care) (1.9.0)	
P(monitor/computer) = marcity + .03( = )=.3966	
plientes vioci ele port oc prilidedes	
"SI as" in all four Cod a VICE "in an ital"	
but the difference between probabilities	
but the difference between probabilities  (or (viv)) vis now smaller.  Za) tf = logiolcountle,d)+1)	
10((n'h)): N13 (2-10)) + 1 (2-10) 2 = (2-201) m) 3	)-
Za) $\xi f = \log_{10} \left( \operatorname{count}(t,d) + 1 \right)$	
P. = (= logio (dF) ) ** ** Calculations done in	
tf-idf= tf x idf. excel sheet.	
65 - (- 5) (6,7.6) xors 1	-
Car 3211 .1551 .3101  ( ) Maracito . NO < ( ) .4610 m/2 .44472	
(2) 3211 .1531 .3101 (2) 4610 (2) .4447214	tage control control we approximate
5 insurance 3148 36686 so 1 1000 "sode" se	paragraphy metro de
pour best 1140 5 1190 5 1190 20 20 1100 0 1910 12 16 13 14	
Twooking Trust words 12 1 was	personal and the second of the second
zustraves du bendess prise de la strodilent	
(cosine $(d_{1}, d_{1}) = 0.6786$	
Cosine $(d_1, d_3) = 0.4401$	energy specific
$cosine(\vec{d}_1,\vec{d}_3) = 0.5509$	)
	C B CONTRACTOR

3a) Simplify CE (y, g) for skipgram model CE(y, g) = - Z. V. log(g.)  $P(O = \omega_0 | \mathbf{I} = \omega_{\mathbf{r}}) = \exp(\alpha \mathbf{i} \omega_0 \cdot \mathbf{i} \omega_{\mathbf{r}})$   $= \exp(\alpha \mathbf{i} \omega_0 \cdot \mathbf{i} \omega_{\mathbf{r}})$   $= (\mathbf{E}(\mathbf{y}, \mathbf{\hat{y}}) = -\log(P(O = \omega_0 | \mathbf{I} = \omega_{\mathbf{r}}))$ 3b) Want to Find DVur [-log(Plo=wol ==wz))] First expand P(0=wolt=vz): JUL [-log/ exp(UTWO.VWI)] De log ( Texp (ut. Vuz)) - Utuo Vwz Now take the derivative we respect to Une of equation 2:

\[
\frac{1}{2} = \frac{\text{Universe}}{\text{Universe}} = \frac{1}{2} - Exp(UTL/VWI) - UTUO WENDERD TEXP(UTL/VWI) = 2 1 2 9 2 - 2 200 = U9 - uu. = 109 - 109



3d) the idea behind negative sampling is to use an loss function which is less computationally expensive than the softmax function which has a summation over the entire vocabulary in the denominator. Negative sampling works similar to binary classification in which the negative samples are used to reduce the probability of positive samples. The value of k in negative sampling is used to determine how many negative samples to use. The loss function for negative sampling is provided below:

log o (Vino VwI) + & Evinpa (w) [log o (-Vini VwI)]

Credit: Mikolov et. al. 2013

Lillill to the total total total debet to the total

6

## Question 4 Written Responses

- a. Done. Submitted to gradescope.
- b. The reason behind the consistent behavior of the first letter in the char ngram model can be attributed to the fact that the first char in the training corpus is F and since we read the entire corpus as one line, F is the only character associated with the starting context (~\*n). The same thing happens for the word level ngram model, since the first word in the training corpus is "In" the model only associates that one word with the starting context (~\*n) and thus "In" always has the highest probability when starting a random text.
- c. Word vs Char on shakespeare sonnets (both with N=2, K=.25, no interpolation)
  - i. Word Model perplexity for shakespeare sonnexts: 60550.778575062905
  - ii. Char Model perplexity for shakespeare sonnexts: 7.9664707724881545
  - iii. From the perplexity values provided above it is evident that when evaluating the performance of the model intrinsically char level models perform better. This being said, when looking at the text generated by each model the extrinsic performance of the word model was better since it produced somewhat coherent text, in the case of the char model the text generated seemed like pure gibberish.
- d. Best Values of perplexity
  - i. Char Ngram (training: Shakespeare\_inputs, dev: Shakespear\_sonnest):5.443296 using N=4, K=.05, and not interpolation
  - ii. Word Ngram (training: Shakespeare\_inputs, dev: Shakespear\_sonnest): 3854.7768 using n=3, K=1.25, and lambda values that are determined by the frequency of the context.
  - iii. Word Ngram (training: train\_e, dev: val\_e): 1496.6144 using N=2, k=.03125, and lambda values of 1 /( N + 1) for all lambdas.
- e. Comparing Char-RNN model and Char Ngram Model
  - i. Best Char Ngram perplexity on shakespeare\_sonnets: 5.443296 using N=4,
     K=.05, and not interpolation
  - ii. Best RNN Ngram perplexity on shakespeare\_sonnets: 6.90747 using hidden\_size=3, epochs=2000, Ir=.005, n\_layers=3
  - iii. Sample Generated by Char-RNN model:

The like the bate will have thew

As a come her have the love the procesters in the saint be my streech.

HARTONE:

A say, my shall but a bent of the part the comes,

What con the prisens our shall thy horders the wand of the done.

QUEE HERNY:

O courter

iv. In terms of intrinsic performance the char ngram model had a better complexity but without a doubt when you compare the results from both models the text generated by the Char-RNN has recognizable words and seems to have some structure to it, unlike the results of the char ngram model.