

## 5. N-grams and Language Models

CS674  
Spring 2007

## Reading Assignment

Jurafsky and Martin Ch. 4  
N-Grams  
from new edition

## Zeros

If  $c(v,w)$  is zero, then  $x_w^v$  is zero.

If the  $vw$  occurs in the testing sample, we get infinite perplexity,  
infinite cross entropy.

Just because we never saw  $vw$  in training, we shouldn't consider it impossible.

## Add one smoothing

Unigram model

$$P(w) =_{\text{def}} c(w) / N$$

$$P'(w) =_{\text{def}} (c(w)+1) / (N + V)$$

$N$  unit counts from the corpus  
 $V$  from add one

## Adjusted counts

$$c^*(w) =_{\text{def}} (c(w) + 1) * N / (N + V)$$

	i	want	to	eat	chinese	food	hunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
hunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Figure 4.1 Bigram counts for eight of the words (out of  $V = 1446$ ) in the Berkeley Restaurant Project corpus of 9332 sentences.

Fig. 4.2 shows the bigram probabilities after normalization (dividing each row

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.0083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Figure 4.2 Bigram probabilities for eight words in the Berkeley Restaurant Project corpus of 9332 sentences.

Now we can compute the probability of sentences like I want English food or I

	i	want	to	eat	chinese	food	lunch	spend
i	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	6	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	1	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

Figure 4.5 Add-one smoothed bigram counts for eight of the words (out of  $V = 1446$ ) in the Berkeley Restaurant Project corpus of 9332 sentences.

For add-one smoothed bigram counts we need to augment the unigram count by

	i	want	to	eat	chinese	food	lunch	spend
i	0.0015	0.21	0.00025	0.0025	0.00025	0.00025	0.00025	0.00075
want	0.0013	0.00042	0.26	0.00084	0.0029	0.0029	0.0025	0.00084
to	0.00078	0.00026	0.0013	0.18	0.00078	0.00026	0.0018	0.055
eat	0.00046	0.00046	0.0014	0.00046	0.0078	0.0014	0.02	0.00046
chinese	0.0012	0.00062	0.00062	0.00062	0.0062	0.052	0.0012	0.00062
food	0.0063	0.00039	0.0063	0.00039	0.00079	0.002	0.00039	0.00039
lunch	0.0017	0.00056	0.00056	0.00056	0.0011	0.00056	0.00056	0.00056
spend	0.0012	0.00058	0.0012	0.00058	0.00058	0.00058	0.00058	0.00058

Figure 4.6 Add-one smoothed bigram probabilities for eight of the words (out of  $V = 1446$ ) in the Berkeley Restaurant Project corpus of 9332 sentences.

It is often convenient to reconstruct the count matrix so we can see how much

	i	want	to	eat	chinese	food	lunch	spend
i	3.8	527	0.64	6.4	0.64	0.64	0.64	1.9
want	1.2	0.39	238	0.78	2.7	2.7	2.3	0.78
to	1.9	0.63	3.1	430	1.9	0.63	4.4	133
eat	0.34	0.34	1	0.34	5.8	1	15	0.34
chinese	0.2	0.098	0.098	0.098	0.098	8.2	0.2	0.098
food	6.9	0.43	6.9	0.43	0.86	2.2	0.43	0.43
lunch	0.57	0.19	0.19	0.19	0.19	0.38	0.19	0.19
spend	0.32	0.16	0.32	0.16	0.16	0.16	0.16	0.16

Figure 4.7 Add-one reconstructed counts for eight of the words (out of  $V = 1446$ ) in the Berkeley Restaurant Project corpus of 9332 sentences.

The above counts in counts and probabilities seem to be much smaller

## Redistributive tax system

- Take a high proportion of income from the rich.
- Redistribute it to the poor as a negative income tax.
- Gives “poor estimates” of the ideal income of high earners.

## Good-Turing smoothing

Count for count

$$N_x = \sum_w 1$$

w such that

$$c(w) = x$$

Number of words that have count c.

## Good-Turing smoothing

$$N_x = \sum 1$$

w such that

$$c(w) = x$$

Any word w which contributes to the sum has an original count  $c(w)=x$ .

Adjust this to  $c^*(w) = (x + 1)N_{x+1} / N_x$

AP Newswire			Berkeley Restaurant		
c (MLE)	$N_x$	$c^*$ (GT)	c (MLE)	$N_x$	$c^*$ (GT)
0	74,671,100,000	0.0000270	0	2,081,496	0.002553
1	2,018,046	0.446	1	5315	0.533960
2	449,721	1.26	2	1419	1.357294
3	188,933	2.24	3	642	2.373832
4	105,668	3.24	4	381	4.081365
5	68,379	4.22	5	311	3.781350
6	48,190	5.19	6	196	4.500000

Figure 4.8 Bigram "frequencies of frequencies" and Good-Turing re-estimations from the 22 million AP bigrams from Church and Gale (1991), and from the Berkeley Restaurant corpus of 9352 sentences.