

Intro to Natural Language Processing

- Last class
 - Unsmoothed n-gram models
- Today:
 - Smoothing
 - » Add-one (Laplacian)
 - » Good-Turing
 - Combining estimators
 - » (Deleted) interpolation
 - » Backoff

Predicting the next word

- Let's go outside and take a ...
 - $P(w_n | w_1^{n-1})$
- Problem?
 - Bigram model
$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-1})$$
 - Trigram model
$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-2} w_{n-1})$$

Training N-gram models

- N-gram models can be trained by counting and normalizing
 - Bigrams
$$P(w_n | w_{n-1}) = \frac{\text{count}(w_{n-1} w_n)}{\text{count}(w_{n-1})}$$
 - General case
$$P(w_n | w_{n-N+1}^{n-1}) = \frac{\text{count}(w_{n-N+1}^{n-1} w_n)}{\text{count}(w_{n-N+1}^{n-1})}$$

Probability of a word sequence

- $P(w_1, w_2, \dots, w_{n-1}, w_n)$
$$P(w_1^n) = P(w_1) P(w_2 | w_1) P(w_3 | w_1^2) \dots P(w_n | w_1^{n-1})$$
$$= \prod_{k=1}^n P(w_k | w_1^{k-1})$$
- Problem?
- Solution: approximate the probability of a word given all the previous words...

N-gram approximations

- Bigram model

$$P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-1})$$

- Trigram model

$$P(w_1^{n-1}) \approx \prod_{k=1}^n P(w_k | w_{k-2} w_{k-1})$$

- N-gram approximation

$$P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-N+1}^{k-1})$$

- Markov assumption: probability of some future event (next word) depends only on a limited history of preceding events (previous words)

Training N-gram models

- N-gram models can be trained by counting and normalizing

- Bigrams

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

- General case

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

- ➔ An example of Maximum Likelihood Estimation (MLE)
 - » Resulting parameter set is one in which the likelihood of the training set T given the model M (i.e. $P(T|M)$) is maximized.

Bigram probabilities

- Problem for the maximum likelihood estimates: sparse data

	I	want	to	eat	Chinese	food	lunch
I	.0023	.32	0	.0038	0	0	0
want	.0025	0	.65	0	.0049	.0066	.0049
to	.00092	0	.0031	.26	.00092	0	.0037
eat	0	0	.0021	0	.020	.0021	.055
Chinese	.0094	0	0	0	0	.56	.0047
food	.013	0	.011	0	0	0	0
lunch	.0087	0	0	0	0	.0022	0

Smoothing

- Need better estimators than MLE for rare events
- Approach
 - Somewhat decrease the probability of previously seen events, so that there is a little bit of probability mass left over for previously unseen events
 - » Smoothing
 - » Discounting methods

Add-one smoothing

- Add one to all of the counts before normalizing into probabilities
- Normal unigram probabilities

$$P(w_x) = \frac{\text{count}(w_x)}{N}$$

- Smoothed unigram probabilities

$$P(w_x) = \frac{\text{count}(w_x) + 1}{N + V}$$

- Adjusted counts

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

Add-one smoothing: bigrams/trigrams

Alternate to adjusted counts

- Adjusted/discounted counts

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

- Discount d_c

$$d_c = \frac{c^*}{c}$$

Add-one bigram counts

- Original counts

	I	want	to	eat	Chinese	food	lunch
I	8	1087	0	13	0	0	0
want	3	0	786	0	6	8	6
to	3	0	10	860	3	0	12
eat	0	0	2	0	19	2	52
Chinese	2	0	0	0	0	120	1
food	19	0	17	0	0	0	0
lunch	4	0	0	0	0	1	0

- New counts

	I	want	to	eat	Chinese	food	lunch
I	9	1088	1	14	1	1	1
want	4	1	787	1	7	9	7
to	4	1	11	861	4	1	13
eat	1	1	3	1	20	3	53
Chinese	3	1	1	1	1	121	2
food	20	1	18	1	1	1	1
lunch	5	1	1	1	1	2	1

Add-one smoothed bigram probabilities

Original

	I	want	to	eat	Chinese	food	lunch
I	.0023	.32	0	.0038	0	0	0
want	.0025	0	.65	0	.0049	.0066	.0049
to	.00092	0	.0031	.26	.00092	0	.0037
eat	0	0	.0021	0	.020	.0021	.055
Chinese	.0094	0	0	0	0	.56	.0047
food	.013	0	.011	0	0	0	0
lunch	.0087	0	0	0	0	.0022	0

Add-one smoothing

	I	want	to	eat	Chinese	food	lunch
I	.0018	.22	.00020	.0028	.00020	.00020	.00020
want	.0014	.00035	.28	.00035	.0025	.0032	.0025
to	.00082	.00021	.0023	.18	.00082	.00021	.0027
eat	.00039	.00039	.0012	.00039	.0078	.0012	.021
Chinese	.0016	.00055	.00055	.00055	.00055	.066	.0011
food	.0064	.00032	.0058	.00032	.00032	.00032	.00032
lunch	.0024	.00048	.00048	.00048	.00048	.00096	.00048

Adjusted bigram counts

Original

	I	want	to	eat	Chinese	food	lunch
I	8	1087	0	13	0	0	0
want	3	0	786	0	6	8	6
to	3	0	10	860	3	0	12
eat	0	0	2	0	19	2	52
Chinese	2	0	0	0	0	120	1
food	19	0	17	0	0	0	0
lunch	4	0	0	0	0	1	0

Adjusted add-one

(#'s are off...)

	I	want	to	eat	Chinese	food	lunch
I	6	740	.68	10	.68	.68	.68
want	2	.42	331	.42	3	4	3
to	3	.69	8	594	3	.69	9
eat	.37	.37	1	.37	7.4	1	20
Chinese	.36	.12	.12	.12	.12	15	.24
food	10	.48	9	.48	.48	.48	.48
lunch	1.1	.22	.22	.22	.22	.44	.22

Too much probability mass is moved

- Estimated bigram frequencies
- AP data, 44million words
- Church and Gale (1991)
- In general, add-one smoothing is a poor method of smoothing
- Much worse than other methods in predicting the actual probability for unseen bigrams

$r = f_{MLE}$	f_{emp}	f_{add-1}
0	0.000027	0.000137
1	0.448	0.000274
2	1.25	0.000411
3	2.24	0.000548
4	3.23	0.000685
5	4.21	0.000822
6	5.23	0.000959
7	6.21	0.00109
8	7.21	0.00123
9	8.26	0.00137