





Add-One Smoothing						
хуа	1	1/3	2	2/29		
xyb	0	0/3	1	1/29		
хус	0	0/3	1	1/29		
xyd	2	2/3	3	3/29		
хуе	0	0/3	1	1/29		
xyz	0	0/3	1	1/29		
Total xy	3	3/3	29	29/29		
600 465 - Intro to NI P	L Fisner			5		

Add-One Smoothing					
300 observations instead of 3 – better data, less smoothing					
хуа	100	100/300	101	101/326	
xyb	0	0/300	1	1/326	
хус	0	0/300	1	1/326	
xyd	200	200/300	201	201/326	
хуе	0	0/300	1	1/326	
xyz	0	0/300	1	1/326	
Total xy	300	300/300	326	326/326	

Add-One Smoothing						
Suppose we're considering 20000 word types, not 26 letters						
1	1/3	2	2/29			
0	0/3	1	1/29			
0	0/3	1	1/29			
2	2/3	3	3/29			
0	0/3	1	1/29			
0	0/3	1	1/29			
3	3/3	29	29/29			
	Smo nsidering 1 0 2 0 2 0 3	Smoothing nsidering 20000 word 1 1/3 0 0/3 0 0/3 2 2/3 0 0/3 0 0/3 0 0/3 3 3/3	Smoothing nsidering 20000 word types, no 1 1/3 2 0 0/3 1 0 0/3 1 2 2/3 3 0 0/3 1 2 2/3 3 0 0/3 1 0 0/3 1 3 3/3 29			

Add-One Smoothing						
As we see m	As we see more word types, smoothed estimates keep falling					
see the abacus	1	1/3	2	2/20003		
see the abbot	0	0/3	1	1/20003		
see the abduct	0	0/3	1	1/20003		
see the above	2	2/3	3	3/20003		
see the Abram	0	0/3	1	1/20003		
see the zygote	0	0/3	1	1/20003		
Total	3	3/3	20003	20003/20003		
600 465 - Intro to NI	PI. Fisner	-	-	8		

Add-Lambda Smoothing

- Suppose we're dealing with a vocab of 20000 words
- As we get more and more training data, we see more and more words that need probability – the probabilities of existing words keep dropping, instead of converging
- This can't be right eventually they drop too low.
- So instead of adding 1 to all counts, add λ = 0.01
- This gives much less probability to those extra events
- But how to pick *best value* for α ? (for the size of our training corpus)
- Try lots of values on a simulated test set! "held-out data"
- Or even better: 10-fold cross validation (aka "jackknifing")
 - Divide data into 10 subsets
 - To evaluate a given alpha:
 - Measure performance on each subset when other 9 are used for training
 Average performance over the 10 subsets tells us how good alpha is



Word type = distinct vocabulary item

Word token = occurrence of that type

Terminology









Good-Turing Smoothing

- Let N_r = # of word types with r training tokens
- Let N = Σ r N_r = total # of training tokens
- Naïve estimate: if x has r tokens, p(x) = ?
 Answer: r/N
- Total naïve probability of all words with r tokens?
 Answer: N_r r / N.
- Good-Turing estimate of this total probability:
 Defined as: N_{r+1} (r+1) / N
- So proportion of novel words in test data is estimated by proportion of singletons in training data.
- Proportion in test data of the N₁ singletons is estimated by
- proportion of the N₂ doubletons in training data. Etc.
- So what is Good-Turing estimate of p(x)?

Use the backoff, Luke! Why are we treating all novel events as the same? p(zygote | see the) vs. p(baby | see the)

- Suppose both trigrams have zero count
- baby beats zygote as a unigram
- the baby beats the zygote as a bigram
- see the baby beats see the zygote ?
- As always for backoff:
- Lower-order probabilities (unigram, bigram) aren't quite what we want
- But we do have enuf data to estimate them & they're better than nothing.

Smoothing + backoff

- Basic smoothing (e.g., add-λ or Good-Turing):
 - Holds out some probability mass for novel events
 - E.g., Good-Turing gives them total mass of N_1/N
 - Divided up evenly among the novel events
- Backoff smoothing
 - Holds out same amount of probability mass for novel events
 - But divide up unevenly in proportion to backoff prob.
 - For p(z | xy):
 - Novel events are types z that were never observed after xy
 - Backoff prob for p(z | xy) is p(z | y) ... which in turn backs off to p(z)!
 - Note: How much mass to hold out for novel events in context xy?
 Depends on whether position following xy is an open class
 - Usually not enough data to tell, though, so aggregate with other contexts (all contexts? similar contexts?)

Deleted Interpolation

- Can do even simpler stuff:
- Estimate p(z | xy) as weighted average of the naïve MLE estimates of p(z | xy), p(z | y), p(z)
- The weights can depend on the context xy
 - If a lot of data are available for the context, then trust p(z | xy) more since well-observed
 - If there are not many singletons in the context, then trust p(z | xy) more since closed-class
 - Learn the weights on held-out data w/ jackknifing