

LANGUAGE MODELING: SMOOTHING

David Kauchak
CS159 – Fall 2024

some slides adapted from
Susan Egan

1

Admin

Assignment 2

- bigram language modeling
- Java
- Can work with partners
 - Anyone looking for a partner?
- 2a: Due Thursday
- 2b: Due next Wednesday
- Style/commenting (JavaDoc)
- Some advice
 - Start now!
 - Spend 1-2 hours working out an example by hand (you can check your answers with me)
 - HashMap

2

Admin

Sumi mentor hours today 6-9pm

Lab next class

Same time. Location TBA (likely upstairs in 229)

3

Today



smoothing
techniques

4

Today

Take home ideas:

- Key idea of smoothing is to redistribute the probability to handle less seen (or never seen) events
 - Still must always maintain a true probability distribution

Lots of ways of smoothing data

Should take into account characteristics of your data!

5

Smoothing

What if our test set contains the following sentence, but one of the trigrams never occurred in our training data?

$P(\text{I think today is a good day to be me}) =$

$P(\text{I} | \langle \text{start} \rangle \langle \text{start} \rangle) \times$
 $P(\text{think} | \langle \text{start} \rangle \text{I}) \times$
 $P(\text{today} | \text{I think}) \times$
 $P(\text{is} | \text{think today}) \times$
 $P(\text{a} | \text{today is}) \times$
 $P(\text{good} | \text{is a}) \times$
 ...

If any of these has never been seen before, prob = 0!

6

Smoothing

$P(\text{I think today is a good day to be me}) =$

$P(\text{I} | \langle \text{start} \rangle \langle \text{start} \rangle) \times$
 $P(\text{think} | \langle \text{start} \rangle \text{I}) \times$
 $P(\text{today} | \text{I think}) \times$
 $P(\text{is} | \text{think today}) \times$
 $P(\text{a} | \text{today is}) \times$
 $P(\text{good} | \text{is a}) \times$
 ...

These probability estimates may be inaccurate. Smoothing can help reduce some of the noise.

7

The general smoothing problem

			modification	probability
see the abacus	1	1/3	?	?
see the abbot	0	0/3	?	?
see the abduct	0	0/3	?	?
see the above	2	2/3	?	?
see the Abram	0	0/3	?	?
...			?	?
see the zygote	0	0/3	?	?
Total	3	3/3	?	?

8

Add-lambda smoothing

A large dictionary makes novel events too probable.

add $\lambda = 0.01$ to all counts

see the abacus	1	1/3	1.01	1.01/203
see the abbot	0	0/3	0.01	0.01/203
see the abduct	0	0/3	0.01	0.01/203
see the above	2	2/3	2.01	2.01/203
see the Abram	0	0/3	0.01	0.01/203
...			0.01	0.01/203
see the zygote	0	0/3	0.01	0.01/203
Total	3	3/3	203	

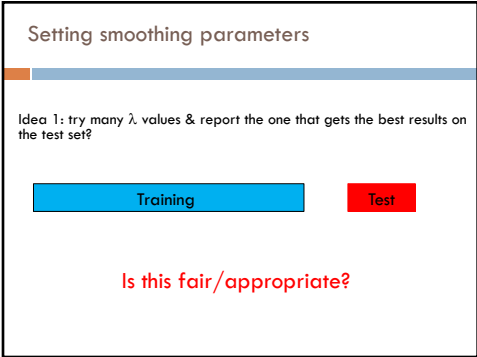
9

Add-lambda smoothing

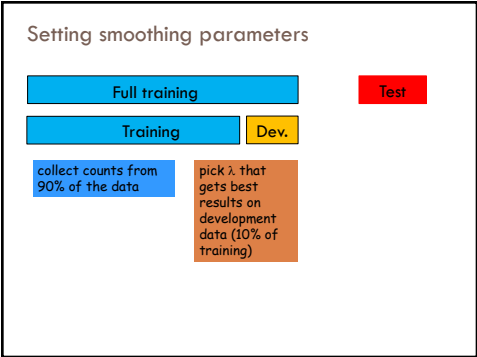
How should we pick lambda?

see the abacus	1	1/3	1.01	1.01/203
see the abbot	0	0/3	0.01	0.01/203
see the abduct	0	0/3	0.01	0.01/203
see the above	2	2/3	2.01	2.01/203
see the Abram	0	0/3	0.01	0.01/203
...			0.01	0.01/203
see the zygote	0	0/3	0.01	0.01/203
Total	3	3/3	203	

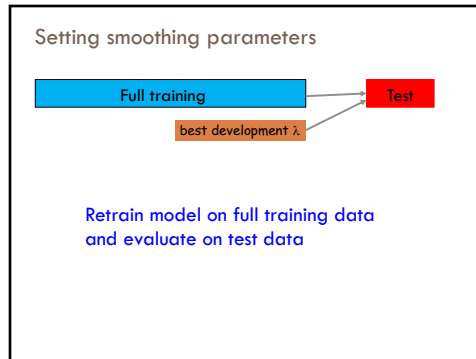
10



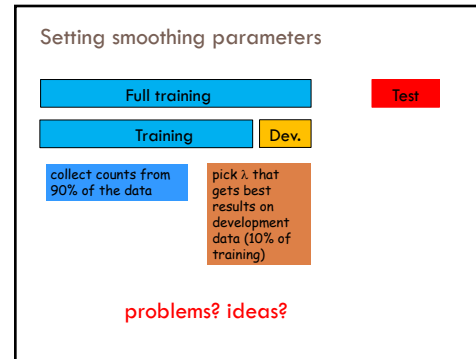
11



12



13



14

Vocabulary

n-gram language modeling assumes we have a fixed vocabulary

□ why?

Probability distributions are over finite events!

see the abacus	1	p(abacus see the)	1/3
see the abbot	0	p(abbot see the)	0/3
see the abduct	0	p(abduct see the)	0/3
see the above	2	p(above see the)	2/3
see the Abram	0	p(Abram see the)	0/3
...			
see the zygote	0	p(zygote see the)	0/3
Total	3		3/3

15

Vocabulary

n-gram language modeling assumes we have a fixed vocabulary

□ why?

Probability distributions are over finite events!

What happens when we encounter a word not in our vocabulary (Out Of Vocabulary)?

- If we don't do anything, prob = 0 (or it's not defined)
- Smoothing doesn't really help us with this!

16

Vocabulary

To make this explicit, smoothing helps us with...

all entries in our vocabulary

see the abacus	1	1.01
see the abbot	0	0.01
see the abduct	0	0.01
see the above	2	2.01
see the Abram	0	0.01
...		0.01
see the zygote	0	0.01

17

Vocabulary

and...

Vocabulary	Counts	Smoothed counts
a	10	10.01
able	1	1.01
about	2	2.01
account	0	0.01
acid	0	0.01
across	3	3.01
...
young	1	1.01
zebra	0	0.01

How can we have words in our vocabulary we've never seen before?

18

Vocabulary

Choosing a vocabulary: **ideas?**

- ▣ Grab a list of English words from somewhere
- ▣ Use all of the words in your training data
- ▣ Use some of the words in your training data
 - for example, all those that occur more than k times

Benefits/drawbacks?

- ▣ Ideally your vocabulary should represent words you're likely to see
- ▣ Too many words: end up washing out your probability estimates (and getting poor estimates)
- ▣ Too few: lots of out of vocabulary

19

Vocabulary

No matter how you chose your vocabulary, you're still going to have out of vocabulary (OOV) words

How can we deal with this?

- ▣ Ignore words we've never seen before
 - Somewhat unsatisfying, though can work depending on the application
 - Probability is then dependent on how many in vocabulary words are seen in a sentence/text
- ▣ Use a special symbol for OOV words and estimate the probability of out of vocabulary

20

Out of vocabulary

Add an extra word in your vocabulary to denote OOV (e.g., <OOV>, <UNK>)

Replace all words in your training corpus not in the vocabulary with <UNK>

- You'll get bigrams, trigrams, etc with <UNK>
 - $p(\langle \text{UNK} \rangle \mid \text{"I am"})$
 - $p(\text{fast} \mid \text{"I } \langle \text{UNK} \rangle \text{"})$

During testing, similarly replace all OOV with <UNK>

21

Choosing a vocabulary

A common approach (and the one we'll use for the assignment):

- Replace the first occurrence of each word by <UNK> in a data set
- Estimate probabilities normally

Vocabulary then is all words that occurred two or more times

This also discounts all word counts by 1 and gives that probability mass to <UNK>

22

Storing the table

How are we storing this table?
Should we store all entries?

see the abacus	1	1/3	1.01	1.01/203
see the abbot	0	0/3	0.01	0.01/203
see the abduct	0	0/3	0.01	0.01/203
see the above	2	2/3	2.01	2.01/203
see the Abram	0	0/3	0.01	0.01/203
...			0.01	0.01/203
see the zygote	0	0/3	0.01	0.01/203
Total	3	3/3	203	

23

Storing the table

Hashtable (e.g. HashMap)

- fast retrieval
- fairly good memory usage

Only store those entries of things we've seen

- for example, we don't store $|V|^3$ trigrams/probabilities

For bigrams we can:

- Store one hashtable with bigrams as keys
- Store a hashtable of hashtables (I'm recommending this)

24

Storing the table: add-lambda smoothing

For those we've seen before:

Unsmoothed (MLE)

$$P(c|ab) = \frac{C(abc)}{C(ab)}$$

add-lambda smoothing

$$P(c|ab) = \frac{C(abc) + \lambda}{C(ab) + ?}$$

see the abacus	1	1/3	1.01	1.01/203
see the abbot	0	0/3	0.01	0.01/203
see the abduct	0	0/3	0.01	0.01/203
see the above	2	2/3	2.01	2.01/203
see the Abram	0	0/3	0.01	0.01/203
...				
see the zygote	0	0/3	0.01	0.01/203
Total	3	3/3	203	

What value do we need here to make sure it stays a probability distribution?

25

Storing the table: add-lambda smoothing

For those we've seen before:

Unsmoothed (MLE)

$$P(c|ab) = \frac{C(abc)}{C(ab)}$$

add-lambda smoothing

$$P(c|ab) = \frac{C(abc) + \lambda}{C(ab) + \lambda V}$$

see the abacus	1	1/3	1.01	1.01/203
see the abbot	0	0/3	0.01	0.01/203
see the abduct	0	0/3	0.01	0.01/203
see the above	2	2/3	2.01	2.01/203
see the Abram	0	0/3	0.01	0.01/203
...				
see the zygote	0	0/3	0.01	0.01/203
Total	3	3/3	203	

For each word in the vocabulary, we pretend we've seen it λ times more (V = vocabulary size).

26

Storing the table: add-lambda smoothing

For those we've seen before:

$$P(c|ab) = \frac{C(abc) + \lambda}{C(ab) + \lambda V}$$

Unseen n-grams: $p(z|ab) = ?$

$$P(z|ab) = \frac{\lambda}{C(ab) + \lambda V}$$

27

Problems with frequency based smoothing

The following bigrams have never been seen:

$p(X | \text{San})$

$p(X | \text{ate})$

Which would add-lambda pick as most likely?

Which would you pick?

28

Witten-Bell Discounting

Some words are more likely to be followed by new words

Diego		food
Francisco		apples
		bananas
San Luis	ate	hamburgers
Jose		a lot
Marcos		for two
		grapes
		...

29

Witten-Bell Discounting

Probability mass is shifted around, depending on the context of words

If $P(w_i | w_{i-1}, \dots, w_{i-m}) = 0$, then the smoothed probability $P_{\text{WBG}}(w_i | w_{i-1}, \dots, w_{i-m})$ is higher if the sequence w_{i-1}, \dots, w_{i-m} occurs with many different words w_k .

30

Problems with frequency based smoothing

The following trigrams have never been seen:

$p(\text{car} | \text{see the })$ $p(\text{zygote} | \text{see the })$

$p(\text{kumquat} | \text{see the })$

Which would add-lambda pick as most likely?
Witten-Bell?

Which would you pick?

34

Better smoothing approaches

Utilize information in lower-order models

trigram	$p(\text{car} \text{see the })$	$p(\text{zygote} \text{see the })$	$p(\text{kumquat} \text{see the })$
bigram	$p(\text{car} \text{the })$	$p(\text{zygote} \text{the })$	$p(\text{kumquat} \text{the })$
unigram	$p(\text{car})$	$p(\text{zygote})$	$p(\text{kumquat})$

35

Better smoothing approaches

Utilize information in lower-order models

Interpolation

- Combine probabilities of lower-order models in some linear combination

Backoff

$$P(z|xy) = \begin{cases} \frac{C^*(xyz)}{C(xy)} & \text{if } C(xyz) > k \\ \alpha(xy)P(z|y) & \text{otherwise} \end{cases}$$

- Often $k = 0$ (or 1)
- Combine the probabilities by "backing off" to lower models only when we don't have enough information

36

Smoothing: simple interpolation

$$P(z|xy) = \lambda \frac{C(xyz)}{C(xy)} + \mu \frac{C(yz)}{C(y)} + (1 - \lambda - \mu) \frac{C(z)}{C(\bullet)}$$

Trigram is very context specific, very noisy

Unigram is context-independent, smooth

Interpolate Trigram, Bigram, Unigram for best combination

How should we determine λ and μ ?

37

Smoothing: finding parameter values

Just like we talked about before, split training data into training and development

Try lots of different values for λ, μ on heldout data, pick best

Two approaches for finding these efficiently

- EM (expectation maximization)
- "Powell search" – see Numerical Recipes in C

38

Backoff models: absolute discounting

$$P_{\text{absolute}}(z|xy) = \begin{cases} \frac{C(xyz) - D}{C(xy)} & \text{if } C(xyz) > 0 \\ \alpha(xy)P_{\text{absolute}}(z|y) & \text{otherwise} \end{cases}$$

Subtract some absolute number from each of the counts (e.g. 0.75)

- How will this affect rare words?
- How will this affect common words?

39

Backoff models: absolute discounting

$$P_{\text{absolute}}(z|xy) = \begin{cases} \frac{C(xy)z - D}{C(xy)} & \text{if } C(xy)z > 0 \\ \alpha(xy)P_{\text{absolute}}(z|y) & \text{otherwise} \end{cases}$$

Subtract some absolute number from each of the counts (e.g. 0.75)

- will have a large effect on low counts (rare words)
- will have a small effect on large counts (common words)

40

Backoff models: absolute discounting

$$P_{\text{absolute}}(z|xy) = \begin{cases} \frac{C(xy)z - D}{C(xy)} & \text{if } C(xy)z > 0 \\ \alpha(xy)P_{\text{absolute}}(z|y) & \text{otherwise} \end{cases}$$

What is $\alpha(xy)$?

41

Backoff models: absolute discounting

trigram model: $p(z|xy)$ (before discounting) trigram model $p(z|xy)$ (after discounting) bigram model $p(z|y)^*$ (*for z where xyz didn't occur)

$$P_{\text{absolute}}(z|xy) = \begin{cases} \frac{C(xy)z - D}{C(xy)} & \text{if } C(xy)z > 0 \\ \alpha(xy)P_{\text{absolute}}(z|y) & \text{otherwise} \end{cases}$$

42

Backoff models: absolute discounting

see the dog	1
see the cat	2
see the banana	4
see the man	1
see the woman	1
see the car	1

$p(\text{cat} | \text{see the}) = ?$
 $p(\text{puppy} | \text{see the}) = ?$

$$P_{\text{absolute}}(z|xy) = \begin{cases} \frac{C(xy)z - D}{C(xy)} & \text{if } C(xy)z > 0 \\ \alpha(xy)P_{\text{absolute}}(z|y) & \text{otherwise} \end{cases}$$

43

Backoff models: absolute discounting

see the dog	1		$p(\text{cat} \mid \text{see the}) = ?$
see the cat	2		
see the banana	4		
see the man	1		
see the woman	1		
see the car	1		

$$\frac{2-D}{10} = \frac{2-0.75}{10} = .125$$

$$P_{\text{absolute}}(z \mid xy) = \begin{cases} \frac{C(xyz)-D}{C(xy)} & \text{if } C(xyz) > 0 \\ \alpha(xy)P_{\text{absolute}}(z \mid y) & \text{otherwise} \end{cases}$$

44

Backoff models: absolute discounting

see the dog	1	0.025	$p(\text{puppy} \mid \text{see the}) = ?$
see the cat	2	0.125	
see the banana	4	0.325	$\alpha(\text{see the}) = ?$
see the man	1	0.025	
see the woman	1	0.025	
see the car	1	0.025	

How much probability mass did we reserve/discount for the bigram model?

$$P_{\text{absolute}}(z \mid xy) = \begin{cases} \frac{C(xyz)-D}{C(xy)} & \text{if } C(xyz) > 0 \\ \alpha(xy)P_{\text{absolute}}(z \mid y) & \text{otherwise} \end{cases}$$

45

Backoff models: absolute discounting

see the dog	1	0.025	$p(\text{puppy} \mid \text{see the}) = ?$
see the cat	2	0.125	
see the banana	4	0.325	$\alpha(\text{see the}) = ?$
see the man	1	0.025	
see the woman	1	0.025	
see the car	1	0.025	

$$\frac{\# \text{ of types starting with "see the"} * D}{\text{count("see the X")}}$$

For each of the unique trigrams, we subtracted D/count("see the") from the probability distribution

$$P_{\text{absolute}}(z \mid xy) = \begin{cases} \frac{C(xyz)-D}{C(xy)} & \text{if } C(xyz) > 0 \\ \alpha(xy)P_{\text{absolute}}(z \mid y) & \text{otherwise} \end{cases}$$

46

Backoff models: absolute discounting

see the dog	1	0.025	$p(\text{puppy} \mid \text{see the}) = ?$
see the cat	2	0.125	
see the banana	4	0.325	$\alpha(\text{see the}) = ?$
see the man	1	0.025	
see the woman	1	0.025	
see the car	1	0.025	

$$\frac{\# \text{ of types starting with "see the"} * D}{\text{count("see the X")}}$$

$$\text{reserved_mass}(\text{see the}) = \frac{6 * D}{10} = \frac{6 * 0.75}{10} = 0.45$$

distribute this probability mass to all bigrams that we are backing off to

$$P_{\text{absolute}}(z \mid xy) = \begin{cases} \frac{C(xyz)-D}{C(xy)} & \text{if } C(xyz) > 0 \\ \alpha(xy)P_{\text{absolute}}(z \mid y) & \text{otherwise} \end{cases}$$

47

Backoff models: absolute discounting

see the dog	1	0.025	$p(\text{puppy} \mid \text{see the}) = ?$
see the cat	2	0.125	$\alpha(\text{see the}) = ?$
see the banana	4	0.325	
see the man	1	0.025	
see the woman	1	0.025	
see the car	1	0.025	

Alternatively, can sum up the discounted probabilities and then see how much probability mass is left

$$P_{\text{absolute}}(z \mid xy) = \begin{cases} \frac{C(xy) - D}{C(xy)} & \text{if } C(xy) > 0 \\ \alpha(xy)P_{\text{absolute}}(z \mid y) & \text{otherwise} \end{cases}$$

48

Backoff models: absolute discounting

see the dog	1	0.025	$p(\text{puppy} \mid \text{see the}) = ?$
see the cat	2	0.125	$\alpha(\text{see the}) = ?$
see the banana	4	0.325	
see the man	1	0.025	
see the woman	1	0.025	
see the car	1	0.025	

Alternatively, can sum up the discounted probabilities and then see how much probability mass is left

sum = 0.55 remaining = 0.45

$$P_{\text{absolute}}(z \mid xy) = \begin{cases} \frac{C(xy) - D}{C(xy)} & \text{if } C(xy) > 0 \\ \alpha(xy)P_{\text{absolute}}(z \mid y) & \text{otherwise} \end{cases}$$

49

Calculating α

We have some number of bigrams we're going to backoff to, i.e. those X where $C(\text{see the } X) = 0$, that is unseen bigrams starting with "see the"

When we backoff, for each of these, we'll be including their probability in the model: $P(X \mid \text{the})$

α is the normalizing constant so that the sum of these probabilities equals the reserved probability mass

$$\alpha(\text{see the})^* \sum_{X: C(\text{see the } X) = 0} p(X \mid \text{the}) = \text{reserved_mass}(\text{see the})$$

50

Calculating α

$$\alpha(\text{see the})^* \sum_{X: C(\text{see the } X) = 0} p(X \mid \text{the}) = \text{reserved_mass}(\text{see the})$$

↓

$$\alpha(\text{see the}) = \frac{\text{reserved_mass}(\text{see the})}{\sum_{X: C(\text{see the } X) = 0} p(X \mid \text{the})}$$

51

Calculating α

We can calculate α two ways

- Based on those we haven't seen:

$$\alpha(\text{see the}) = \frac{\text{reserved_mass}(\text{see the})}{\sum_{X:C(\text{see the } X) = 0} p(X|\text{the})}$$

- Or, more often, based on those we do see:

$$\alpha(\text{see the}) = \frac{\text{reserved_mass}(\text{see the})}{1 - \sum_{X:C(\text{see the } X) > 0} p(X|\text{the})}$$

52

Calculating α in general: trigrams

$p(X | A B)$

Calculate the reserved mass

$$\text{reserved_mass}(\text{bigram: } A B) = \frac{\text{\# of types starting with } A B * D}{\text{count}(A B X)}$$

where X is any "word"

Calculate the sum of the backed off probability. For bigram "A B":

$$1 - \sum_{X:C(A B X) > 0} p(X|B) \quad \text{either is fine, in practice, the left is easier} \quad \sum_{X:C(A B X) = 0} p(X|B)$$

Calculate α

$$\alpha(A B) = \frac{\text{reserved_mass}(A B)}{1 - \sum_{X:C(A B X) > 0} p(X|B)}$$

1 - the sum of the bigram probabilities of those trigrams that we saw starting with bigram A B

53

Calculating α in general: bigrams

$p(X | A)$

Calculate the reserved mass

$$\text{reserved_mass}(\text{unigram: } A) = \frac{\text{\# of types starting with } A * D}{\text{count}(A X)}$$

where X is any "word"

Calculate the sum of the backed off probability. For unigram "A":

$$1 - \sum_{X:C(A X) > 0} p(X) \quad \text{either is fine in practice, the left is easier} \quad \sum_{X:C(A X) = 0} p(X)$$

Calculate α

$$\alpha(A) = \frac{\text{reserved_mass}(A)}{1 - \sum_{X:C(A X) > 0} p(X)}$$

1 - the sum of the unigram probabilities of those bigrams that we saw starting with word A

54

Calculating backoff models in practice

Store the α s in another table

- If it's a trigram backed off to a bigram, it's a table keyed by the bigrams
- If it's a bigram backed off to a unigram, it's a table keyed by the unigrams

Compute the α s during training

- After calculating all of the probabilities of seen unigrams/bigrams/trigrams
- Go back through and calculate the α s (you should have all of the information you need)

During testing, it should then be easy to apply the backoff model with the α s pre-calculated

55

Backoff models: absolute discounting

the Dow Jones	10	$p(\text{jumped} \mid \text{the Dow}) = ?$
the Dow rose	5	What is the reserved mass?
the Dow fell	5	

$$\frac{\# \text{ of types starting with "the Dow" * D}}{\text{count("the Dow")}}$$

$$\text{reserved_mass}(\text{the Dow}) = \frac{3 * D}{20} = \frac{3 * 0.75}{20} = 0.115$$

$$\alpha(\text{the Dow}) = \frac{\text{reserved_mass}(\text{see the})}{1 - \sum_{X: C(\text{the Dow } X) > 0} p(X \mid \text{the})}$$

56

Backoff models: absolute discounting

$$\text{reserved_mass} = \frac{\# \text{ of types starting with bigram} * D}{\text{count}(\text{bigram})}$$

Two nice attributes:

- ▣ decreases if we've seen more bigrams
 - should be more confident that the unseen trigram is no good
- ▣ increases if the bigram tends to be followed by lots of other words
 - will be more likely to see an unseen trigram

57