# Language Modeling

CMSC 473/673 Spring 2017
Bryan Wilkinson

# Discussion of Kilgarriff Paper

- Does anyone disagree with Kilgarriff?
- What are the alternatives?
  - UkWaC and WaCkypedia
  - UMBC WebBase

# Language Modeling

- A Language Model is a collection of probabilities of some unit of language
  - Word and character level models are the most common.
- Why do we need to know the probability of language?
  - Which spelling is more likely?
  - Which verb is more likely?
- What is the probability of language?
  - Of a word?
  - Of a sentence?

# Counting Words

- Before we can calculate probability, we need to count things.
  - Seems like a simple task but,
- What is a word?
  - Still an open question, but mostly agreed upon
  - Is *a lot* one word or two?
  - How about *post office*?
- English is easy to count
  - What about a language with more complex morphology (We will talk more about this later)
  - [Iyewičhamačheča](#) (Lakota "I resemble you" or "you resemble me")
- Other minor decisions
  - Does capitalization matter?
  - How about plurals? Or other endings?

- Counting is the hard part

$$P(w) = \frac{C(w)}{N}$$

# In Class Activity

# Markov Assumption

- The future doesn't depend on the past
  - But the future **does** depend on the past, so we will just limit the amount of past we look at
- In NLP this means when calculating the probability of a word given a context, we often just look at the preceding one or two words. ( or none!)
- In terms of probability this looks like
  - P(*homework | the dog ate my*) ≈ P(*homework*)
  - P(*homework | the dog ate my*) ≈ P(*homework | my*)
  - P(*homework | the dog ate my*) ≈ P(*homework | ate my*)
- We call these small chunks of words N-grams
  - Unigram = homework
  - Bigram = my homework
  - Trigram = ate my homework
- What are all the bigrams and trigrams in the sentence "the dog ate my homework"

- Calculating probability of an N-Gram is almost as easy as a word
- Technically we are estimating using Maximum Likelihood Estimation

$$P(w|w_{n-1}) = \frac{C(w_{n-1}w_n)}{\sum_w C(w_{n-1}w)}$$

$$P(w|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

- To calculate the probability of a sentence, we use the chain rule

$$P(w_1 w_2 w_3 ... w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)...P(w_n|w_1, w_2, w_3...w_{n-1})$$

- P(Today is the second day of class) = P(Today | <start>) * P(is | Today <start>) * P(the | is Today <start>) * P(second | the is Today <start>) * P(day | second the is Today <start>) * P(of | day second the is Today <start>) * P(class | of day second the is Today <start>) * P(<end> | class of day second the is Today <start>)

- To calculate the probability of a sentence, we use the chain rule

$$P(w_1 w_2 w_3 ... w_n) = P(w_1)P(w_2|w_1)P(w_3|w_2)...P(w_n|w_{n-1})$$

- P(Today is the second day of class) ≈ P(Today | <start>) * P(is | Today) * P(the | is) * P(second | the ) * P(day | second ) * P(of | day ) * P(class | of ) * P(<end> | class )

# Probability of a Sentence

- The following counts are from Google Books N-Gram Corpus v2
  - Via corpus.byu.edu

| C(<start>,Today) | |
|---|---|
| C(Today, is) | 49215 |
| C(is, the) | 101249451 |
| C(the, second) | 17216997 |
| C(second, day) | 396983 |
| C(day, of) | 5639617 |
| C(of, class) | 672495 |
| C(class,<end>) | 2703083 |

| C(<start>) | |
|---|---|
| C(Today) | 3251311 |
| C(is) | 1368855691 |
| C(the) | 7726878625 |
| C(second) | 42738935 |
| C(day) | 78211120 |
| C(of) | 5035745089 |
| C(class) | 27260152 |
| C(<end>) | 7805955213 |

# Probability of a Sentence

| | |
|---|---|
| P(Today \| <start> ) | |
| P(is \| Today) | |
| P(the \| is) | |
| P(second \| the ) | |
| P(day \| second ) | |
| P(of \| day ) | |
| P(class \| of ) | |
| P(<end> \| class ) | |

- What additional data would we need to calculate the probability of :
  - Today is the eighth day of class ?
  - Today is the 57th day of class ?
  - Today is the somethingth day of class?

# Evaluating Language Models

- LMs give us probabilities of a sentence
  - How do we know if they are working well
  - What is the "correct" probability of a given sentence
- We could evaluate them using some task that has LMs as part of it
  - Test a machine translation, speech recognition, etc system and see which LM does better
  - Can conflates any errors
- So we try to evaluate LMs as a standalone system
  - This is hard
  - We do this using perplexity, a measure from information theory.

# Perplexity

- Wikipedia (and probably statistics books) define perplexity as 2 raised to the entropy of the distribution
  - Technically correct, but not very intuitive (to me at least)
- Jurafsky gives the definition as the probability an LM gives to a test set, normalized by the number of words in that test set.
  - The idea is that words and sequences of words that actually exist in text should have a relatively high probability
- The assumption is that perplexity should correlate to performance on the end task
  - Often is the case!
  - But not always!

$$\sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i|w_1...w_{i-1})}}$$

- This data comes from a system I am working on that has a perplexity of about 3

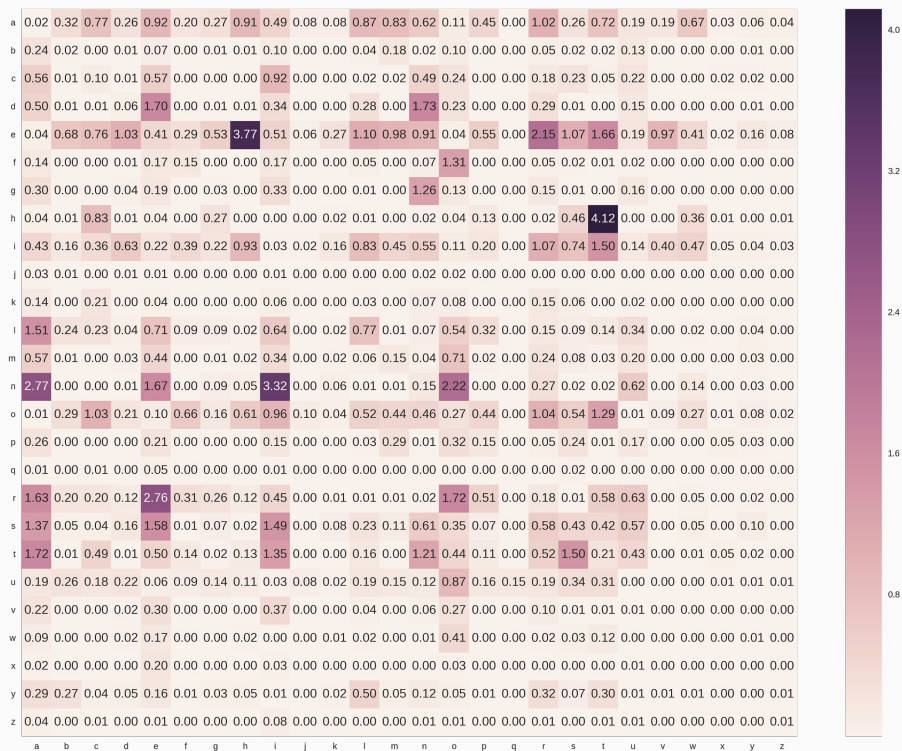| Input | Output |
|---|---|
| Who are you ? | Who are you ? |
| This caused problem like the formation and growth of slums which most of the time is not safe due to their unhealthy environment . | This caused problem like the both and growth of of which most of the time not not not be to their . . |
| Tokyo is often referred to as a city, but is officially known and governed as a "metropolitan prefecture" . | New is often referred to as a city , but is officially known and workforce as a " the prefecture " . |
| The Tokyo metropolitan government administers the 23 Special Wards of Tokyo | The financially of government classified the 00 Special Wards of New |

# Character Level Language Models

- What are they useful for?
  - Language ID
  - Spell Check
- They are smaller
  - Just need to count letters and a few other symbols
- Becoming more popular with Neural Networks
  - [The Unreasonable Effectiveness of Recurrent Neural Networks](#)

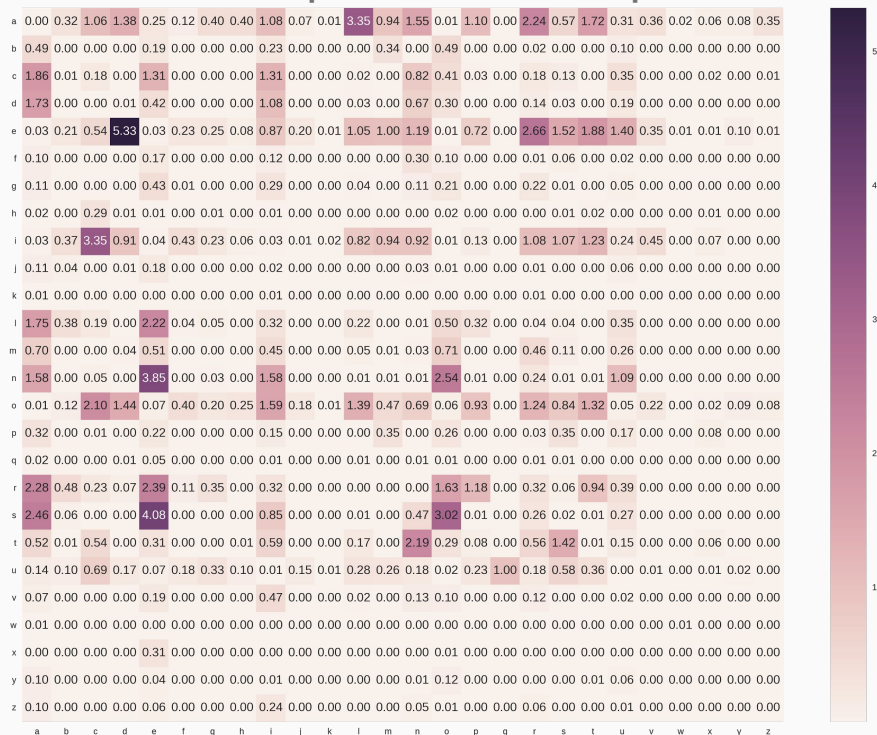- Do they represent actual linguistic units?

# Character LM Example

- Let's say we want to build a simple language ID system
- Bigram model over 2GB of wackypeadia produces the following counts

# Character LM Example

- Let's say we want to build a simple language ID system
- Bigram model over 2GB of French UN text produces the following counts

# Character LM Example

- Let's say we want to build a simple language ID system
- Bigram model over 2GB of Spanish UN text produces the following counts

- Let's try the following sentences and see the results
  - Today is the second day of class
  - La universidad tienen 50 años.
  - Marius se sentit fier de cet inconnu.
  - Melissa Villaseñor is on SNL

- What happens when we put in a language we don't have a model for
  - What can this tell us?

# Reminder

- HW 0 "Due" Tuesday