

# Language Modeling

# Generalization and zeros



# The Shannon Visualization Method

- Choose a random bigram  
( $\langle s \rangle$ ,  $w$ ) according to its probability
- Now choose a random bigram  
( $w$ ,  $x$ ) according to its probability
- And so on until we choose  $\langle /s \rangle$
- Then string the words together

$\langle s \rangle$  I  
 I want  
 want to  
 to eat  
 eat Chinese  
 Chinese food  
 food  $\langle /s \rangle$   
 I want to eat Chinese food



# Approximating Shakespeare

## Unigram

To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have  
 Every enter now severally so, let  
 Hill he late speaks; or! a more to leg less first you enter  
 Are where exeunt and sighs have rise excellency took of.. Sleep knave we. near; vile like

## Bigram

What means, sir. I confess she? then all sorts, he is trim, captain.  
 Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.  
 What we, hath got so she that I rest and sent to scold and nature bankrupt, nor the first gentleman?

## Trigram

Sweet prince, Falstaff shall die. Harry of Monmouth's grave.  
 This shall forbid it should be branded, if renown made it empty.  
 Indeed the duke; and had a very good friend.  
 Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.

## Quadrigram

King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;  
 Will you not tell me who I am?  
 It cannot be but so.  
 Indeed the short and the long. Marry, 'tis a noble Lepidus.



## Shakespeare as corpus

- $N=884,647$  tokens,  $V=29,066$
- Shakespeare produced 300,000 bigram types out of  $V^2= 844$  million possible bigrams.
  - So 99.96% of the possible bigrams were never seen (have zero entries in the table)
- Quadrigrams worse: What's coming out looks like Shakespeare because it *is* Shakespeare



# The wall street journal is not shakespeare (no offense)

## Unigram

Months the my and issue of year foreign new exchange's september were recession ex-  
change new endorsed a acquire to six executives

## Bigram

Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor  
would seem to complete the major central planners one point five percent of U. S. E. has  
already old M. X. corporation of living on information such as more frequently fishing to  
keep her

## Trigram

They also point to ninety nine point six billion dollars from two hundred four oh six three  
percent of the rates of interest stores as Mexico and Brazil on market conditions



# The perils of overfitting

- N-grams only work well for word prediction if the test corpus looks like the training corpus
  - In real life, it often doesn't
  - We need to train robust models that generalize!
  - One kind of generalization: Zeros!
    - Things that don't ever occur in the training set
      - But occur in the test set



# Zeros

- Training set:
  - ... denied the allegations
  - ... denied the reports
  - ... denied the claims
  - ... denied the request
- Test set
  - ... denied the offer
  - ... denied the loan

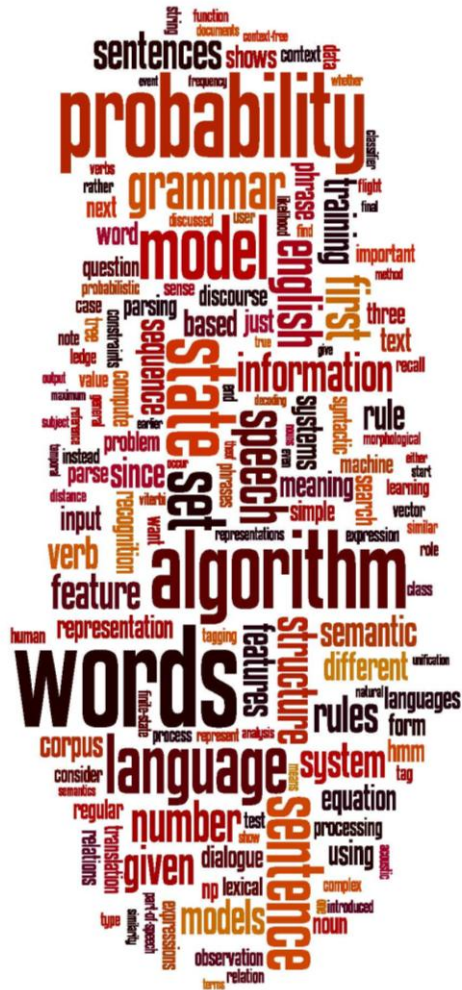
$$P(\text{"offer"} \mid \text{denied the}) = 0$$



# Zero probability bigrams

- Bigrams with zero probability
  - mean that we will assign 0 probability to the test set!
- And hence we cannot compute perplexity (can't divide by 0)!





# Language Modeling

Generalization and  
zeros