

# CMPT 413

## Computational Linguistics

Anoop Sarkar

<http://www.cs.sfu.ca/~anoop>

## $n$ -grams

- A simple model of language
- Computes a probability for observed input
- Probability is likelihood of observation being generated by the same source as the training data
- Such a model is often called a *language model*

# An example

- Let's consider an example we've seen before: *spelling correction*

*... was called a “stellar and versatile **acress** whose combination of sass and glamour has defined her ...*

KCG model best guess is **acres**

# An example

- A language model can take the context into account:
  - ... was called a “stellar and versatile **acress** whose combination of sass and glamour has defined her ...*
  - ... was called a “stellar and versatile **acres** whose combination of sass and glamour has defined her ...*
  - ... was called a “stellar and versatile **actress** whose combination of sass and glamour has defined her ...*
- Each sentence is a sequence  $w_1, \dots, w_n$ . Task is to find  $P(w_1, \dots, w_n)$ .

# Another example

physical Brainpower not plant is chief , now a 's asset , . firm  
, a Brainpower not now chief asset firm 's is . plant physical ,  
chief a physical , . firm not , Brainpower plant is asset 's now  
not plant Brainpower now physical 's . a chief , asset firm , is  
plant Brainpower is now , , not . firm a 's physical asset chief  
physical is 's plant firm not chief . Brainpower now asset , , a  
Brainpower , not physical plant , is now a firm 's chief asset .

Each sentence is a sequence  $w_1, \dots, w_n$ .

Task is to find  $P(w_1, \dots, w_n)$ .

How can we compute  $P(w_1, \dots, w_n)$

- Apply the *Chain rule*
- $P(w_1, \dots, w_n) = P(w_1) \cdot P(w_2 \mid w_1) \cdot P(w_3 \mid w_1, w_2) \dots P(w_n \mid w_1, \dots, w_{n-1})$
- Each of these probabilities can be estimated (using frequency counts) from *training data*
- **But** we need to apply these probabilities on unseen *test data*
- The curse of dimensionality: **sparse data**

# The Markov Assumption

*a stellar and versatile **acres** whose combination of*

*$P(a) \cdot P(\text{stellar} \mid a) \cdot P(\text{and} \mid a, \text{stellar}) \cdot$*

*$P(\text{versatile} \mid a, \text{stellar}, \text{and}) \cdot$*

*$P(\text{acres} \mid a, \text{stellar}, \text{and}, \text{versatile}) \cdot$*

*$P(\text{whose} \mid a, \text{stellar}, \text{and}, \text{versatile}, \text{acres}) \dots$*

*a stellar and versatile **acres** whose combination of*

*$P(a) \cdot P(\text{stellar} \mid a) \cdot P(\text{and} \mid a, \text{stellar}) \cdot P(\text{versatile} \mid \text{stellar}, \text{and}) \cdot$*

*$P(\text{acres} \mid \text{and}, \text{versatile}) \cdot P(\text{whose} \mid \text{versatile}, \text{acres}) \dots$*

## $n$ -grams

- 0th order Markov model:  $P(w_i)$  called a *unigram* model
- 1st order Markov model:  $P(w_i \mid w_{i-1})$  called a *bigram* model
- 2nd order Markov model:  $P(w_i \mid w_{i-2}, w_{i-1})$  called a *trigram* model



# Parameter size

**Corpus:** <s> said the joker to the thief

N (tokens) = 7      |V| = 6

$$\begin{aligned} p(joker|the) &= \frac{p(the, joker)}{p(the)} \\ &= \frac{\frac{f(the, joker)}{\text{num of bigrams}}}{\frac{f(the)}{\text{num of unigrams}}} = \frac{f(the, joker)}{f(the)} \end{aligned}$$

## $n$ -grams

- How many possible distinct probabilities will be needed?, i.e. **parameter values**
- Total number of **word tokens** in our training data
- Total number of unique words: **word types** is our vocabulary size

# $n$ -gram Parameter Sizes

- Let  $V$  be the vocabulary, size of  $V$  is  $|V|$
- $P(W_i = x)$ , how many different values for  $W_i$ 
  - $|V| = 3 \times 10^3$
- $P(W_i = x \mid W_j = y)$ , how many different values for  $W_i, W_j$ 
  - $|V|^2 = 9 \times 10^6$
- $P(W_i = x \mid W_k = z, W_j = y)$ , how many different values for  $W_i, W_j, W_k$ 
  - $|V|^3 = 27 \times 10^9$

# Parameter size

**Corpus:**      <s> said the joker to the thief

$$|V| = 6$$

**Bigrams:**    max num of parameters =  $|V|^2 = 36$

said | <s>  
the | said  
joker | the  
to | joker  
the | to  
thief | the

$$\text{observed} = W_T = 6 \quad \ll 36$$

# $n$ -gram model of Jane Austen

- Three novels by Jane Austen: *Emma*, *Sense and Sensibility*, *Pride and Prejudice*
- Removed punctuation and kept paragraph structure
- Trained a trigram model on this text



# $n$ -gram model of Jane Austen

$f(3\text{gram})$	$f(2\text{gram})$	$f(1\text{gram})$	$w_0$	$w_1$	$w_2$
378	518	10381	I	do	not
366	1366	10381	I	am	sure
214	1917	9182	in	the	world
202	572	6917	she	could	not
189	462	2751	would	have	been
174	184	10381	I	dare	say
173	179	5758	as	soon	as
173	357	11135	a	great	deal
171	332	7573	it	would	be
155	945	3017	could	not	be

# *n*-gram model of Jane Austen

3gram $\frac{f(w_0, w_1, w_3)}{f(w_0, w_1)}$	2gram $\frac{f(w_0, w_1)}{f(w_0)}$	1gram $\frac{f(w_0)}{N}$	$w_0$	$w_1$	$w_2$
0.72	0.04	0.016	I	do	not
0.26	0.13	0.016	I	am	sure
0.11	0.20	0.014	in	the	world
0.35	0.08	0.011	she	could	not
0.40	0.16	0.004	would	have	been
0.94	0.01	0.016	I	dare	say
0.96	0.03	0.009	as	soon	as
0.48	0.03	0.018	a	great	deal
0.51	0.04	0.012	it	would	be
0.16	0.31	0.004	could	not	be

# $n$ -gram model of Jane Austen

$f(3\text{gram})$	$f(2\text{gram})$	$f(1\text{gram})$	$w_0$	$w_1$	$w_2$
1	1	1	favor	of	your
1	1	1	peerage	his	wealth
1	1	1	stagnation	Mrs	Elton's
1	1	1	genteelly	and	paid
1	1	1	adept	in	the
1	1	1	deckers	now	in
1	1	1	oracle	Fanny's	explanations
1	1	1	Ashworth	is	too
1	1	1	puddles	with	impatient
1	1	1	Harringtons	to	come
1	1	1	roasted	No	coffee
1	1	1	coherent	Dearest	Lizzy



# $n$ -gram model of Jane Austen

3gram	2gram	1gram	$w_0$	$w_1$	$w_2$
0.00039	0.13	0.029	of	the	Middletons
0.00039	0.13	0.029	of	the	Meryton
0.00039	0.13	0.029	of	the	Lucases
0.00039	0.13	0.029	of	the	London
0.00039	0.13	0.029	of	the	Irish
0.00039	0.13	0.029	of	the	History
0.00039	0.13	0.029	of	the	First
0.00039	0.13	0.029	of	the	Esquire
0.00039	0.13	0.029	of	the	Elegant
0.00039	0.13	0.029	of	the	Dashwoods
0.00039	0.13	0.029	of	the	Crown

# Summary

- $n$ -grams define a probability model over sequences
  - we have seen examples of sequences of words, but you can define  $n$ -grams over sequences of characters or other sequences
- $n$ -grams deal with sparse data by using the Markov assumption
- The number of parameters increase rapidly when the value of  $n$  is increased for  $n$ -grams but the data cannot keep up with the parameter size.