### Improving Text Simplification Language Modeling Using Unsimplified Text Data

David Kauchak Computer Science Department Middlebury College dkauchak@middlebury.edu

### **Text simplification**



Any intelligent fool can make things bigger, more complex, and more violent. It takes a touch of genius and a lot of courage to move in the opposite direction.



Simpler is better.

### Simpler is better



### Goal:

Reduce the reading complexity of text by incorporating more accessible vocabulary and structure while maintaining the content.

I find forest colored chicken ovum and smoked pork thigh to be dietarily disturbing.



model



I find forest colored chicken ovum and smoked pork thigh to be dietarily disturbing.



I do not like green eggs and ham.





I find forest colored chicken ovum and smoked pork thigh to be dietarily disturbing.

translation model

language model

length model

model



I do not like green eggs and ham.



I find forest colored chicken ovum and smoked pork thigh to be dietarily disturbing.



model

### **Data availability**



## How much data is available to train a *simple* English language model?

~0.5 millions sentences



### **Data availability**



# How much data is available to train an English language model?

A lot more.





simple n-grams found in normal Wikipedia

n-grams	simple <b>→</b> normal % overlap
1	96%
2	80%
3	68%
4	61%
5	55%

Sentence aligned corpus (137K sentence pairs)





Sentence aligned corpus (137K sentence pairs)

simple n-grams found in normal Wikipedia

n-grams	simple <b>→</b> normal % overlap
1	96%
2	80%
3	68%
4	61%
5	55%

#### Good news:

- some reasonable overlap
- It's English



96%

80%

68%

**61%** 

55%



simple n-grams found in normal Wikipedia

Sentence aligned corpus (137K sentence pairs)

Possibly bad news: a lot of missing data!



WIKIPEDIA The Free Encyclopedia	Simple English WIKIPEDIA			
(normal)	(simple)	n-grams	simple <b>-</b> ≯normal	normal <b>-</b> ≯simple
		1	96%	87%
		2	80%	68%
		3	68%	58%
		4	61%	51%
		5	55%	46%

Sentence aligned corpus (137K sentence pairs)

Bad news: different distributions over English!





How do these distribution differences affect language modeling performance?

Is unsimplified data useful for simple language modeling?

What is the best way to utilize unsimplified data?

## **Document Aligned Corpus**



#### en.wikipedia.org/wiki/England

#### simple.wikipedia.org/wiki/England

Simple English WIKIPEDIA





**simple-only**: simple English Wikipedia sentences

normal-only: English Wikipedia sentences

**simple-X+normal**: X simple sentences combined with varying amounts of normal sentences





23% improvement in perplexity by adding normal data to simple data





## normal data helps even more if the simple data is limited



### Language model adaptation

Linearly interpolated language model:

$$p_{\text{interpolated}}(W_i \mid W_{i-2}W_{i-1}) = \lambda p_{\text{normal}}(W_i \mid W_{i-2}W_{i-1}) + (1-\lambda)p_{\text{simple}}(W_i \mid W_{i-2}W_{i-1})$$

normal-only

simple-only





~24% improvement in perplexity over models trained with ALL available simple data by using normal data

### **Task 2: Lexical simplification**



SemEval 2012 task:

With the physical market as *tight* as it has been in memory, silver could fly at any time.

#### **Candidates**

constricted pressurised low high-strung tight



Task: ranker

Human simplicity ranking
tight

low constricted pressurised high-strung

### **Task 2: Lexical simplification**

With the physical market as *constricted* as it has been ... With the physical market as *pressurised* as it has been ... With the physical market as *low* as it has been ... With the physical market as *high-strung* as it has been ...

With the physical market as *tight* as it has been ...







### **Task 2: Evaluation**







### **Lexical simplification results**





### Less simple data





number of additional normal sentences



23% improvement over simple-only model!

### Why does normal data help?



Our guess: more *n*-grams

How many more *n*-grams are seen in normal data compared to simple?



### Why does normal data help?



How many more *n*-grams are seen in normal data compared to simple?

Perplexity test data

Lexical simplification data

re

normal contains:

I

unigrams	9.4% more	6.2% more
bigrams	24% more	56% more
trigrams	46% more	117% more

### **Application matters**



Optimal  $\lambda$  (weighting between simple and normal) for linearly interpolated models:

Perplexity task

 $\lambda = 0.5$ 

An equal balance between simple and normal models Lexical simplification task

 $\lambda = 0.98$ 

A very strong bias towards the simple model





Unsimplified data **is** useful for simple English language modeling

>23% improvement on both perplexity and lexical simplification tasks over model using *ALL simple data* available

LM domain adaption techniques are important, but are application specific

Data available:

http://www.cs.middlebury.edu/~dkauchak/simplification/

### **Open questions**



How much unsimplified data can we utilize?

How does source/domain affect perfomance?

How does the LM quality affect other simplification applications (e.g. full sentence simplification)?

Better LM domain adaptation techniques.