A Word-Complexity Lexicon and A Neural Readability Ranking Model for Lexical Simplification

Mounica Maddela and Wei Xu









Department of Computer Science and Engineering INPUT: Applesauce is a **puree** made of apples.

OUTPUT: Applesauce is a <u>soft paste</u>. It is made of apples.

Text Simplification

INPUT: Applesauce is a **puree** made of apples.

OUTPUT: Applesauce is a <u>soft paste</u>. It is made of apples.



Applications

- Reading assistance for children, non-native speakers and disabled.
- Improve other NLP tasks (MT, summarization ...)

INPUT: Applesauce is a puree made of apples.

OUTPUT: Applesauce is a soft paste. It is made of apples.

INPUT: Applesauce is a **puree** made of apples.

OUTPUT: Applesauce is a soft paste. It is made of apples.

Complex Word Identification

INPUT: Applesauce is a **puree** made of apples.

OUTPUT: Applesauce is a **soft paste.** It is made of apples. **liquidized sauce thick liquid**

Complex Word Identification - Substitution Generation

INPUT: Applesauce is a **puree** made of apples.

OUTPUT: Applesauce is a <u>soft paste</u>. It is made of apples. thick liquid liquidized sauce

Complex Word Identification - Substitution Generation - Substitution Ranking

A Large Word-complexity Lexicon

• 15,000 English words w/ human ratings



• predict relative complexity for any given words or phrases



A Pairwise Neural Ranking Model

• improve the state-of-the-art significantly for all lexical simplification tasks



Complex Word Identification - Substitution Generation - Substitution Ranking

(% is relative error reduction)

Previous Work

Rely on **heuristics and corpus level features** to measure word complexity

• Word length

(Shardlow 2013, Biran et. al. 2011, and many others)

• Word frequency in corpus

(Bott et. al. 2011, Kajiwara et. al. 2013, Horn et. al. 2014, and many others)

• Language model probability

(Glavas & Stajner 2015, Paetzold & Special 2016/17, and many others)

Weakness of Previous Work



- pundit > professional
 - alien > stranger

* based on 2272 lexical paraphrases sampled from PPDB

Weakness of Previous Work



- folly > foolishness
- scheme > outline
- distress > discomfort

* based on 2272 lexical paraphrases sampled from PPDB

A Large Word-complexity Lexicon

- 15,000 most frequent English words from Google 1T ngram corpus
- Rated on a 6-point Likert scale



- 15,000 most frequent English words from Google 1T ngram corpus
- Rated on a 6-point Likert scale



- 11 annotators (non-native speakers)
- ► 5 ~ 7 ratings for each word
- 2.5 hours to rate 1000 words





- Inter-annotator agreement is 0.64 (Pearson correlation)
- One annotator rating vs. mean of the rest

Word	Score	A1	A2	A3	A 4	A 5
muscles	1.6	2	1	2	2	1
pattern	2.4	2	3	1	1	3
educational	3.2	3	3	3	3	4
cortex	4.2	4	4	4	4	5
assay	5.8	6	6	6	5	6

difference (one vs. rest)

- < 0.5 for 47% of annotations
- < 1.0 for 78% of annotations
- < 1.5 for 93% of annotations

Evaluation* - Complex Word Identification

- Complex Word Identification Shared Task BEA@NAACL'18
- 34879 sentences from Wikipedia and news articles
- 27299 training, 3328 development, 4252 test instances

Input	The whale was sensing him with sound pulses .
Output	[Complex, simple]

* see paper for full evaluation on 3 lexical simplification tasks and 5 benchmark datasets

- Complex Word Identification Shared Task 2018
- 27299 training, 3328 development, 4252 test instances

	F-score	Accuracy
Senses	62.3	54.1
SimpleWiki Frequency	63.3	61.6
Length	65.9	67.7
(Yimam et al. 2017)	66.6	76.7
(Paetzold et al. 2016)	73.8	78.7

- Complex Word Identification Shared Task 2018
- 27299 training, 3328 development, 4252 test instances

F-score	Accuracy
62.3	54.1
63.3	61.6
65.9	67.7
66.6 +	1.6 76.7 +2
73.8	78.7
67.5	69.8
-	F-score 62.3 63.3 65.9 66.6 73.8 67.5

- Complex Word Identification Shared Task 2018
- 27299 training, 3328 development, 4252 test instances



* statistically significant (p < 0.01) based on the paired bootstrap test

A Pairwise Neural Ranking Model



Word-Complexity Lexicon Score 0/1 binary indicator

word length word frequency number of syllables ngram probabilities



Word-Complexity Lexicon Score 0/1 binary indicator

word length word frequency number of syllables ngram probabilities

PPDB paraphrase score word2vec cosine similarity











0.91 = P(more_complex | **adversary** - **enemy**)

 $P > 0 \Rightarrow w_a$ is more complex than w_b

 $P < 0 \implies w_a$ is simpler than w_b

P indicates complexity difference

 $\langle w_a : adversary , w_b : enemy \rangle$

Neural Readability Ranking Model



- English Lexical Simplification Shared Task SemEval 2012
- 300 training sentences, 1710 test sentences

Input	There were also pieces that would have been <u>terrible</u> in any environment.
(Paetzold & Specia 2017)	awful, very bad, dreadful
Our Model + Our Lexicon	very bad, awful, dreadful
Gold truth	very bad, awful, dreadful

** see paper for full evaluation on 3 lexical simplification tasks and 5 benchmark datasets

- English Lexical Simplification Shared Task SemEval 2012
- 300 training sentences, 1710 test sentences

		Precision@1	Pearson
heuristics	(Biran et al. 2011)	51.3	0.505
SVM	(Jauhar & Specia 2012)	60.2	0.575
heuristics	(Kajiwara et al. 2013)	60.4	0.649
SVM	(Horn et al. 2014)	63.9	0.673
heuristics	(Glavaš & Štajner 2015)	63.2	0.644
SVM	(Paetzold & Specia 2015)	65.3 🔨 +0.	2 0.677 - +0.002
neural	(Paetzold & Specia 2017)	65.6	0.679
)+1.	7)+0.03
neural Our	Model + Lexicon + Gaussian	67.3*/	0.714*/

* statistically significant (p < 0.05) based on the paired bootstrap test

- English Lexical Simplification Shared Task SemEval 2012
- 300 training sentences, 1710 test sentences

		Precision@1	Pearson
heuris	tics (Biran et al. 2011)	51.3	0.505
SVM	(Jauhar & Specia 2012)	60.2	0.575
heuris	tics (Kajiwara et al. 2013)	60.4	0.649
SVM	(Horn et al. 2014)	63.9	0.673
heuris	tics (Glavaš & Štajner 2015)	63.2	0.644
SVM	(Paetzold & Specia 2015)	65.3 🔨 +0	.2 0.677 🔨 +0.00
neural	(Paetzold & Specia 2017)	65.6	0.679
			7
neural	Our Model + Gaussian	66.6	$(0.702^{*})^{+0.03}$
neural	Our Model + Lexicon + Gaussian	67.3*/	0.714*/

* statistically significant (p < 0.05) based on the paired bootstrap test

- English Lexical Simplification Shared Task SemEval 2012
- 300 training sentences, 1710 test sentences

		Precision@1	Pearson
heuris	tics (Biran et al. 2011)	51.3	0.505
SVM	(Jauhar & Specia 2012)	60.2	0.575
heuris	tics (Kajiwara et al. 2013)	60.4	0.649
SVM	(Horn et al. 2014)	63.9	0.673
heuris	tics (Glavaš & Štajner 2015)	63.2	0.644
SVM	(Paetzold & Specia 2015)	65.3 🔨 +0	.2 0.677 - +0.00
neural	(Paetzold & Specia 2017)	65.6	0.679
neural	Our Model	65.4	0.682
neural	Our Model + Gaussian	66.6	$(0.702^{*})^{+0.03}$
neural	Our Model + Lexicon + Gaussian	67.3*/	0.714*/

* statistically significant (p < 0.05) based on the paired bootstrap test

Evaluation - Error Analysis

Input	The colonies of one <u>strain</u> appeared smooth.
(Paetzold & Specia 2017)	sort, type, breed, variety
Our Model + Our Lexicon	type, sort, breed, variety
Gold truth	type, sort, variety, breed

Input	No damage or <u>casualties</u> were	reported.
(Paetzold & Specia 2017)	injuries, accidents, deaths,	fatalities
Our Model + Our Lexicon	injuries, deaths, accidents,	fatalities
Gold truth	deaths, injuries, accidents,	fatalities

SimplePPDB++

• 14.1 million paraphrase rules w/ improved complexity ranking scores

Paraphrase Rule		
\rightarrow self-supporting	0.93	
self-reliant \rightarrow self-sufficient	0.48	
→ self-sustainable complex	-0.60	
→ possible	0.94	
viable \rightarrow realistic	0.15	
→ plausible	-0.91	
→ in-depth review	0.89	
detailed assessement \rightarrow careful examination	0.28	
→ comprehensive evaluation	-0.87	

Thanks

• Word-Complexity Lexicon & SimplePPDB++ are available!

day	1.0	MIN 1 (simple)
convenient	2.4	
transmitted	3.2	
cohort	4.3	
assay	5.8	MAX 6 (complex)

- PyTorch Code for the Neural Ranking model is also available! https://github.com/mounicam/lexical_simplification
- Contacts: Mounica Maddela & Wei Xu (Ohio State University)

A Word-Complexity Lexicon and A Neural Readability Ranking Model for Lexical Simplification



t-SNE visualization of the complexity scores, ranging between 1.0 and 6.0

Word-Complexity Lexicon

Coverage over Penn Treebank (~1.1 million words)



Single feature value : f(w) = 0.41, $f(w) \in [0,1]$

Vectorized feature : f(w) = [-0.0, 0.44, 0.54, -0.02, -0.0]



Single feature value : f(w) = 0.41, $f(w) \in [0,1]$

Vectorized feature : f(w) = [-0.0, 0.44, 0.54, -0.02, -0.0]



Single feature value : f(w) = 0.41, $f(w) \in [0,1]$

Vectorized feature : f(w) = [-0.0],



Single feature value : f(w) = 0.41, $f(w) \in [0,1]$

Vectorized feature : $f(w) = [\sim 0.0, 0.44]$



Single feature value : f(w) = 0.41, $f(w) \in [0,1]$

Vectorized feature : f(w) = [~0.0, 0.44, 0.54]



Single feature value : f(w) = 0.41, $f(w) \in [0,1]$

Vectorized feature : f(w) = [-0.0, 0.44, 0.54, -0.02, -0.0]



Substitution Ranking - Correct Examples

Our Model predicts the correct output

Input	The <u>concept</u> of a "picture element" dates to the earliest days of television.
(Paetzold & Specia 2017)	theory, thought, idea
Our Model + Our Lexicon	idea, thought, theory
Gold truth	idea, thought, theory

Our Model handles phrases better than previous SOTA.

Input	There were also pieces that would have been <u>terrible</u> in any environment.
(Paetzold & Specia 2017)	awful, very bad, dreadful
Our Model + Our Lexicon	very bad, awful, dreadful
Gold truth	very bad, awful, dreadful