Fully Delexicalized Contexts for Syntax-Based Word Embeddings

Jenna Kanerva¹, Sampo Pyysalo² and Filip Ginter¹

¹Dept of IT - University of Turku, Finland ²Lang. Tech. Lab - University of Cambridge

turkunlp.github.io

Abstract

- We propose fully delexicalized contexts derived from syntactic trees to train word embeddings
- We demonstrate and evaluate our embeddings compared to vanilla word2vec
 - Nearest neighbours
 - Correlation to human judgement
 - Dependency parsing

Outline

- Related work
- Our word embedding contexts
- Motivation
- Experiments and evaluation

Related work

- Word2vec Mikolov et al. (2013)
- Word2vecf Levy and Goldberg (2014)
 Syntactic context

Vanilla word2vec (Mikolov et al. 2013)

Hope you enjoyed your weekend in Houston

Output layer: Hope, you, your, weekend

Hidden layer

Input layer: enjoyed

Hope you enjoyed your weekend in Houston

* Output layer size: vocabulary

Word2vecf (Levy & Goldberg 2014)



Output layer: Hope/ccomp^{gov}, you/nsubj, weekend/obj

* Output layer size: ~vocabulary

Hidden layer

Input layer: enjoyed



6

Our method



Output layer: ccomp^{gov}, nsubj, obj, VERB, Mood=Ind, Tense=Past, VerbForm=Fin

Hidden layer

Input layer: enjoyed

* Output layer size: pos tags + features + dependency types



Motivation

1) How much semantics can be learnt without the actual words?

2) Does task-specific training help?

3) Unified treebank annotations → Universal/multilingual word embeddings?

Experiments

- Training our word embeddings for 45 languages
- Inspecting nearest neighbours
- Comparing our embeddings to vanilla word2vec
 - Correlation with human judgement
 - Dependency parsing (closely related task)

Data

- <u>Word embedding training:</u> Automatically parsed raw text collection (Ginter et al. 2017)
- <u>Parser training + evaluation:</u> Universal
 Dependencies v2.0 treebanks (Nivre et al. 2017)
- <u>Word similarity evaluation</u>: evaluation service of 13 human judgement datasets (Faruqui and Dyer 2014)

Evaluation: Nearest neighbours

france	jesus	xbox	reddish	scratched	megabits	
belgium	christ	playstation	brownish	knicked	megabit	
luxembourg	jesus.	ps3	yellowish	bruised	kilobits	
nantes	god	ps4	greenish	nicked	gigabits	
marseille	ahnsahnghong	xbox360	pinkish	scuffed	mbps	vanilla
bretagne	jesuschrist	wii	grayish	chewed	mbits	word2vec
boulogne	y'shua	xbla	bluish	sandpapere	d terabits	
poitou	christ	psvita	-orange	scratches	mbit	
rouen	christ.	titanfall	orangish	brusied	kbits	
paris	jesus	xboxone	greyish	scraped	kilobit	
toulouse	yeshua	gamecube	mid-brown	thwacked	megabytes	
france	iocus	vbov	reddish	scratched	megahits	
	Jesus	AUUA	reduisii	seruteneu	meguons	
lebanon	osama	vbox	greenish	snatched	megabytes	
lebanon australia	osama napoleon	vbox whitesox	greenish grayish	snatched touched	megabytes microseconds	
lebanon australia england	osama napoleon ophelia	vbox whitesox matchbox	greenish grayish bluish	snatched touched punched	megabytes microseconds hectares	
lebanon australia england bolivia	osama napoleon ophelia gautama	vbox whitesox matchbox firefox	greenish grayish bluish greyish	snatched touched punched deflected	megabytes microseconds hectares tonnes	Our
lebanon australia england bolivia scotland	osama napoleon ophelia gautama scipio	vbox whitesox matchbox firefox wmp	greenish grayish bluish greyish pinkish	snatched touched punched deflected warmed	megabytes microseconds hectares tonnes microns	Our delexicalized
lebanon australia england bolivia scotland estonia	osama napoleon ophelia gautama scipio sauron	vbox whitesox matchbox firefox wmp audiovox	greenish grayish bluish greyish pinkish yellowish	snatched touched punched deflected warmed levelled	megabytes microseconds hectares tonnes microns micrograms	Our delexicalized
lebanon australia england bolivia scotland estonia switzerland	osama napoleon ophelia gautama scipio sauron chandragupta	vbox whitesox matchbox firefox wmp audiovox virtualbox	greenish grayish bluish greyish pinkish yellowish brownish	snatched touched punched deflected warmed levelled booted	megabytes microseconds hectares tonnes microns micrograms litres	Our delexicalized vectors
lebanon australia england bolivia scotland estonia switzerland finland	osama napoleon ophelia gautama scipio sauron chandragupta claudius	vbox whitesox matchbox firefox wmp audiovox virtualbox equinox	greenish grayish bluish greyish pinkish yellowish brownish blackish	snatched touched punched deflected warmed levelled booted stalked	megabytes microseconds hectares tonnes microns micrograms litres megawatts	Our delexicalized vectors
lebanon australia england bolivia scotland estonia switzerland finland slovenia	osama napoleon ophelia gautama scipio sauron chandragupta claudius jamarcus	vbox whitesox matchbox firefox wmp audiovox virtualbox equinox rotax	greenish grayish bluish greyish pinkish yellowish brownish blackish temperate	snatched touched punched deflected warmed levelled booted stalked ditched	megabytes microseconds hectares tonnes microns micrograms litres megawatts gallons	Our delexicalized vectors
lebanon australia england bolivia scotland estonia switzerland finland slovenia algeria	osama napoleon ophelia gautama scipio sauron chandragupta claudius jamarcus olivia	vbox whitesox matchbox firefox wmp audiovox virtualbox equinox rotax hmp	greenish grayish bluish greyish pinkish yellowish brownish blackish temperate redish	snatched touched punched deflected warmed levelled booted stalked ditched swallowed	megabytes microseconds hectares tonnes microns micrograms litres megawatts gallons bushels	Our delexicalized vectors

Evaluation: Nearest neighbours

france	jesus	xbox	reddish	scratched	megabits	
belgium	christ	playstation	brownish	knicked	megabit	
luxembourg	jesus.	ps3	yellowish	bruised	kilobits	
nantes	god	ps4	greenish	nicked	gigabits	
marseille	ahnsahnghong	xbox360	pinkish	scuffed	mbps	vanilla
bretagne	jesuschrist	wii	grayish	chewed	mbits	word2vec
boulogne	y'shua	xbla	bluish	sandpapered	d terabits	
poitou	christ	psvita	-orange	scratches	mbit	
rouen	christ.	titanfall	orangish	brusied	kbits	
paris	jesus	xboxone	greyish	scraped	kilobit	
toulouse	yeshua	gamecube	mid-brown	thwacked	megabytes	
france	jesus	xbox	reddish	scratched	megabits	
lebanon	osama	vbox	greenish	snatched	megabytes	
australia	napoleon	whitesox	grayish	touched	microseconds	
england	ophelia	matchbox	bluish	punched	hectares	Our
bolivia	gautama	firefox	greyish	deflected	tonnes	Our
scotland	scipio	wmp	pinkish	warmed	microns	delexicalized
estonia	sauron	audiovox	yellowish	levelled	micrograms	Vectore
switzerland	chandragupta	virtualbox	brownish	booted	litres	vectors
finland	claudius	equinox	blackish	stalked	megawatts	
slovenia	iamarcus	rotax	temperate	ditched	gallons	
	J				0	
algeria	olivia	hmp	redish	swallowed	bushels	

Evaluation: human judgement

	Correlation		Pairs			
Dataset	word2vec	Our vectors	Found	Total	Reference	
WordSim-353	0.7083	0.2350	353	353	Finkelstein et al. (2001)	
WordSim-353-SIM	0.7677	0.4033	203	203	Agirre et al. (2009)	
WordSim-353-REL	0.6691	0.1318	252	252	Agirre et al. (2009)	
MC-30	0.7028	0.2929	30	30	Miller and Charles (1991)	
RG-65	0.6801	0.0593	65	65	Rubenstein and Goodenough (1965)	
Rare-Word	0.4250	0.1998	2006	2034	Luong et al. (2013)	
MEN	0.7397	0.2027	3000	3000	Bruni et al. (2012)	
MTurk-287	0.6958	0.3474	287	287	Radinsky et al. (2011)	
MTurk-771	0.6406	0.1336	771	771	Halawi et al. (2012)	
YP-130	0.3882	0.0464	130	130	Yang and Powers (2006)	
SimLex-999	0.3376	0.1004	999	999	Hill et al. (2016)	
Verb-143	0.3633	0.2425	144	143	Baker et al. (2014)	
SimVerb-3500	0.2175	0.0476	3500	3500	Gerz et al. (2016)	

 Our vectors not as good as word2vec but correlation is still positive

language	baseline	word2vec	diff to baseline	syntax-based	diff to baseline
Ancient_Greek	56.61	57.93	+1.32	58.18	+1.57
Ancient_Greek-PROIEL	72.35	72.48	+0.13	72.67	+0.32
Arabic	72.88	73.91	+1.03	74.00	+1.12
Basque	69.02	69.74	+0.72	69.93	+0.91
Bulgarian	83.90	84.29	+0.39	85.18	+1.28
Catalan	85.15	85.01	-0.14	85.31	+0.16
Chinese	68.48	68.83	+0.35	69.06	+0.58
Croatian	76.08	75.98	-0.10	77.35	+1.27
Czech-CAC	8375	83 58	-0.17	84 54	+0.79
Czech-CLTT	69.58	68.92	-0.66	72.19	+2.61
Czech	84.47	84 24	.0.23	84 69	+0.22
Danish	75.18	74.63	-0.55	74.99	_0.19
Dutch-LassySmall	75.67	75.01	-0.66	76.68	+1.01
Dutch	74.73	75 21	+0.48	75.00	+0.27
English	79.66	80.20	+0.40	80.64	+0.27
English LinES	74.62	74.35	0.27	75.50	+0.93
English DarTUT	75.72	75.21	-0.27	76.20	+0.97
English-rai i O i	60.65	61.90	-0.51	62.20	+0.40
Estonian	00.05	01.89	+1.24	0.3.22	+2.57
Finnish Eiserich ETD	75.70	15.19	+0.09	11.35	+1.05
Finnish-FTB	70.42	/0.08	+0.20	11.12	+1.50
French	86.08	85.71	-0.37	86.5.5	+0.45
French-Sequoia	82.30	82.58	+0.28	82.65	+0.35
Galician	77.58	77.34	-0.24	78.21	+0.63
German	73.10	73.12	+0.02	72.87	-0.23
Greek	79.04	77.93	-1.11	79.93	+0.89
Hebrew	76.88	77.38	+0.50	78.52	+1.64
Hindi	87.09	86.82	-0.27	87.38	+0.29
Hungarian	65.59	66.40	+0.81	68.44	+2.85
Indonesian	74.39	72.84	-1.55	73.59	-0.80
Italian	85.44	84.98	-0.46	84.96	-0.48
Italian-ParTUT	78.21	78.74	+0.53	79.92	+1.71
Japanese	93.09	93.09	+0.00	93.23	+0.14
Korean	56.42	62.70	+6.28	63.72	+7.30
Latin-ITTB	71.15	71.72	+0.57	72.98	+1.83
Latin-PROIEL	70.08	69.76	-0.32	69.89	-0.19
Latvian	64.01	64.56	+0.55	66.16	+2.15
Norwegian-Bokmaal	83.91	83.44	-0.47	84.18	+0.27
Norwegian-Nynorsk	82.32	81.65	-0.67	81.89	-0.43
Old_Church_Slavonic	73.56	71.22	-2.34	71.40	-2.16
Persian	80.38	79.56	-0.82	80.86	+0.48
Polish	79.42	80.62	+1.20	81.21	+1.79
Portuguese-BR	85.55	86.11	+0.56	86.26	+0.71
Portuguese	83.64	84.49	+0.85	84.93	+1.29
Romanian	79.82	79.77	-0.05	80.30	+0.48
Russian	75.41	76.00	+0.59	77.48	+2.07
Russian-SynTagRus	8676	86.58	-0.18	87.71	+0.95
Slovak	75 30	75.65	+0.26	76.55	+1.16
Slownian	80.62	80.87	+0.25	81.39	+0.76
Spanish-AnCora	84 17	84 55	+0.38	84 31	+0.14
Spanish	84.24	92.95	0.40	84.11	0.22
Swadich LinES	74.35	74 72	+0.37	75.34	+0.00
Swedish-Lines	72.20	74.72	+0.57	73.54	+0.99
Turdai ah	56.00	56.24	+0.80	14.75	+1.30
Unda	76.00	76.22	+0.24	26.26	+1.75
Vistaman	10.98	10.23	-0.75	10.20	-0.72
vietnamese	22.82	50.20	+0.41	55.22	-0.63
Average	-	-	+0.16	-	+0.88

UDPipe parser with three different pre-trained word embeddings

- <u>Baseline</u>: word2vec
 trained on treebank data
- word2vec trained on raw text collection
- <u>Our</u> trained on automatically analysed raw text collection

language	baseline	word2vec	diff to baseline	syntax-based	diff to baseline
Ancient_Greek	56.61	57.93	+1.32	58.18	+1.57
Ancient_Greek-PROIEL	72.35	72.48	+0.13	72.67	+0.32
Arabic	72.88	73.91	+1.03	74.00	+1.12
Basque	69.02	69.74	+0.72	69.93	+0.91
Bulgarian	83.90	84.29	+0.39	85.18	+1.28
Catalan	85.15	85.01	-0.14	85.31	+0.16
Chinese	68.48	68.83	+0.35	69.06	+0.58
Croatian	76.08	75.98	-0.10	77.35	+1.27
Czech-CAC	83.75	83.58	-0.17	84 54	+0.79
Czech-CLTT	69.58	68 92	-0.66	72.19	+2.61
Czech	84.47	84 24	0.23	84.69	+0.22
Danish	75.18	74.63	-0.55	74.09	-0.19
Dutch J assySmall	75.67	75.01	-0.66	76.68	+1.01
Dutch	74.73	75.21	+0.48	75.00	+0.27
English	79.66	80.20	+0.40	80.64	+0.27
English English Lin EC	79.00	24.25	0.07	26.60	+0.90
English DerTUT	74.02	74.35	-0.27	75.59	+0.97
English-Pari Oi	15.12	(1.00	-0.51	/0.20	+0.48
Estonian	60.65	01.89	+1.24	6.3.22	+2.57
Finnish	75.70	15.19	+0.09	11.35	+1.65
Finnish-FIB	/6.42	/0.68	+0.26	11.12	+1.50
French	86.08	85.71	-0.37	86.53	+0.45
French-Sequoia	82.30	82.58	+0.28	82.65	+0.35
Galician	77.58	77.34	-0.24	78.21	+0.63
German	73.10	73.12	+0.02	72.87	-0.23
Greek	79.04	77.93	-1.11	79.93	+0.89
Hebrew	76.88	77.38	+0.50	78.52	+1.64
Hindi	87.09	86.82	-0.27	87.38	+0.29
Hungarian	65.59	66.40	+0.81	68.44	+2.85
Indonesian	74.39	72.84	-1.55	73.59	-0.80
Italian	85.44	84.98	-0.46	84.96	-0.48
Italian-ParTUT	78.21	78.74	+0.53	79.92	+1.71
Japanese	93.09	93.09	+0.00	93.23	+0.14
Korean	56.42	62.70	+6.28	63.72	+7.30
Latin-ITTB	71.15	71.72	+0.57	72.98	+1.83
Latin-PROIEL	70.08	69.76	-0.32	69.89	-0.19
Latvian	64.01	64.56	+0.55	66.16	+2.15
Norwegian-Bokmaal	83.91	83.44	-0.47	84.18	+0.27
Norwegian-Nynorsk	82.32	81.65	-0.67	81.89	-0.43
Old_Church_Slavonic	73.56	71.22	-2.34	71.40	-2.16
Persian	80.38	79.56	-0.82	80.86	+0.48
Polish	79.42	80.62	+1.20	81.21	+1.79
Portuguese-BR	85.55	86.11	+0.56	86.26	+0.71
Portuguese	83.64	84.49	+0.85	84.93	+1.29
Romanian	79.82	79.77	-0.05	80.30	+0.48
Russian	75.41	76.00	+0.59	77.48	+2.07
Russian-SynTagRus	86.76	86.58	-0.18	87.71	+0.95
Slovak	75 30	75.65	+0.26	76.55	+1.16
Slovenian	80.62	80.87	+0.25	81.38	+0.76
Spanish-AnCora	84 17	84 55	+0.38	84 31	+0.14
Spanish	84 34	83.95	-0.40	84.11	.0.22
Swadich LinES	74.35	74 72	+0.37	75.34	+0.00
Sweatsh-Lines	72.20	74.72	+0.57	73.34	+0.99
Turchish	15.59	14.25	+0.80	14.15	+1.50
Turkish Unda	76.00	30.24	+0.24	31.13	+1.75
N	70.98	10.23	-0.75	70.20	-0.72
vietnamese	22.82	50.20	+0.41	55.22	-0.63
Average	-	-	+0.10	-	+0 XX

Green if our is better than word2vec, and difference to baseline is positive

word2vec +0.16 better than baseline on average
 o diff to baseline between -1.55% and +6.28%
 o 31 treebanks positive, and 23 negative

Our embeddings +0.88 better than baseline
 o diff to baseline between -0.80% and +7.30%
 o 45 treebanks positive, and 9 negative

- pre-trained embeddings does not automatically increase parsing performance across languages
- delexicalized syntactic embeddings lead to higher performance as well as generalize better across languages when evaluated in closely related task

Parser accuracy vs. quality of the embeddings

- Our word embeddings are trained on automatically parsed data
 - How does the baseline parser accuracy affect the quality of the word embeddings?
- Bootstrapping:
 - Baseline parser → parse raw text → embeddings → better parser → parse raw text → new embeddings → even better parser?

Bootstrapping on Finnish

	Baseline	Iteration 1	Iteration 2
Finnish	75.70	77.35	77.57

Small improvement with the second iteration model

→ UDPipe not optimal parser for this study as POS tags and morphological features are not revised

 Baseline UDPipe not competitive with state-of-the-art on Finnish
 75.7% compared to 83-84%

• What if we use raw data parsed with (near) state-of-the-art parser?

- Raw data: Finnish Internet Parsebank
 ~3.6 billion token collection of web crawled data
- Finnish-dep-parser
 - Omorfi rule-based morphological analyzer
 - Marmot tagger
 - Mate-tools graph-based dependency parser
 - UD v1.2
 - LAS estimated to be ~82%

Warning! Numbers not comparable!

- Numbers are not comparable to our main result table!
 - Different version of UD (v1.2 compared to v2.0)
 - Raw text collection more than three times bigger

• UDPipe + our embeddings trained on Finnish Internet Parsebank: 82.21%

• UDPipe + word2vec embeddings trained on Finnish Internet Parsebank: 78.35%

• UDPipe baseline: ~76.5%

Conclusions

 Fully delexicalized context for word embedding training

• Bit surprisingly, these embeddings are able to capture also semantic aspects

 Improve parsing accuracy and generalize better than standard word2vec embeddings

Thanks!

turkunlp.github.io

Academy of Finland Kone Foundation University of Turku Graduate School