





A New General Deep Learning Approach for Natural Language Processing

Hui Jiang

Department of Electrical Engineering and Computer Science York University, Toronto, Canada

Joint work with Shiliang Zhang (USTC), Quan Liu (USTC), Mingbin Xu (York), Joseph Sanu (York)

Outline

- Introduction
 - Deep Learning for NLP
- A New General Approach: FOFE-net
 - FOFE: lossless encoding
 - DNNs: universal approximator
- NLP Applications
 - Word Embedding
 - Language Modelling
 - Named Entity Recognition and Mention Detection in KBP-EDL
 - Coreference Resolution in Winograd
- On-going work
 - Entity Linking
 - Syntactic/Semantic Parsing
- Conclusions





the more the better

compact representative

neural networks





Word: word embedding

sentence/paragraph/document:
variable-length word sequences



Deep Learning for NLP: Difficulties

- How to handle variable-length sequences?
 - Extracting fixed-size features —> information loss
 - Using CNN/RNN/LSTM as sequence labelling
 —> Not flexible in use; slow and hard to learn
- How to deal with *data sparsity* problem?
 - *Maybe* more training data + powerful models



- Introduction
 - Deep Learning for NLP
- A New General Approach: FOFE-net
 - FOFE: lossless encoding
 - DNNs: universal approximator
- NLP Applications
 - Word Embedding
 - Language Modelling
 - Named Entity Recognition and Mention Detection in KBP-EDL
 - Coreference Resolution in Winograd
- On-going work
 - Entity Linking
 - Syntactic/Semantic Parsing
- Conclusions



Fixed-size Ordinally-Forgetting Encoding (FOFE)

- FOFE: encode any variable-length sequence into a fixed-size representation
 - Unique and invertible —> no information loss (theoretically guaranteed)
 - Elegant and automatic —> no feature engineering

WORD	1-OF-K		
mone	101 K	ANY SEQUENCE	FOFE
w_0	1000000	W ₆	0,0,0,0,0,0,1
w_1	0100000	W_6, W_A	$0, 0, 0, 0, 1, 0, \alpha$
<i>w</i> ₂	0010000	W_{c}, W_{4}, W_{F}	$0.0.0.0.\alpha.1.\alpha^2$
<i>W</i> ₃	0001000	W_{c} W_{t} W_{r} W_{o}	$1000 \alpha^2 \alpha \alpha^3$
W_4	0000100	W, W, W, W, W,	$\alpha = 0 = 0 = 0 = \alpha^3 = 1 + \alpha^2 = \alpha^4$
W_5	0000010	<i>w</i> ₆ , <i>w</i> ₄ , <i>w</i> ₅ , <i>w</i> ₀ , <i>w</i> ₅	$a_{1}^{2}, 0, 0, 0, 0, 1 + a_{1}^{4}, a_{2}^{4} + a_{3}^{3}, a_{5}^{5}$
We	0000001	$w_6, w_4, w_5, w_0, w_5, w_4$	$\alpha^{-}, 0, 0, 0, 1 + \alpha^{-}, \alpha + \alpha^{-}, \alpha^{-}$

 $\mathbf{z}_t = \alpha \cdot \mathbf{z}_{t-1} + \mathbf{e}_t \ (1 \le t \le n)$

FOFE: theoretical guarantee

• FOFE is *almost unique* for any *a* (0<*a*<1):

Theorem 1 If the forgetting factor α satisfies $0 < \alpha \leq 0.5$, FOFE is unique for any K and T.

Theorem 2 For $0.5 < \alpha < 1$, given any finite values of K and T, FOFE is almost unique everywhere for $\alpha \in (0.5, 1.0)$, except only a finite set of countable choices of α .

A New General Deep Learning Approach for Natural Language Processing

- FOFE-based lossless encoding
 - any text input —> a fixed-size FOFE vector
- Deep neural networks (DNN):
 - universal approximator (Cybenko, '89): learned to map to any NLP targets
- Labelled training data:
 - human labelling; semi-supervised; distant supervision

FOFE-net for NLP



- Theoretically sound
- No feature engineering
- General methodology
 - (almost) all NLP tasks
- Enough labelled data
 - human performance (?)

Applicable NLP tasks

- Word Embedding
- Chunking
 - Tokenization
 - Segmentation
- Word Sense Disambiguation
- POS Tagging
- Named Entity Recognition
- Semantic Role Labelling

- Syntactic Parsing
- Semantic Parsing
- Language Modelling
- Text Categorization
- Coreference Resolution
- Information Extraction
- Relationship Extraction
- more ...

Applicable NLP tasks

- Word Embedding
- Chunking
 - Tokenization
 - Segmentation
- Word Sense Disambiguation
- POS Tagging
- Named Entity Recognition
- Semantic Role Labelling

- Syntactic Parsing
- Semantic Parsing
- Language Modelling
- Text Categorization
- Coreference Resolution
- Relationship Extraction
- Information Extraction
- more ...

Applicable NLP tasks

- Word Embedding
- Chunking

•

- Tokenization
 - Segmentation
- Word Sense Disambiguation
- POS Tagging
- Named Entity Recognition
- Semantic Role Labelling

- Syntactic Parsing
- Semantic Parsing
- Language Modelling
- Text Categorization
- Coreference Resolution
- Relationship Extraction
- Information Extraction
- more ...

Outline

- Introduction
 - Deep Learning for NLP
- A New General Approach: FOFE-net
 - FOFE: lossless encoding
 - DNNs: universal approximator
- NLP Applications
 - Word Embedding
 - Language Modelling
 - Named Entity Recognition and Mention Detection in KBP-EDL

14

- Coreference Resolution in Winograd
- On-going work
 - Entity Linking
 - Syntactic/Semantic Parsing
- Conclusions





FOFE for Word Embedding

 Distributional hypothesis (Harris, 1954): words that appear in similar contexts have similar meaning.

• Using FOFE to model contexts:



FOFE for Word Embedding

Generate word-context (the averaged FOFE) matrix (one word per line):



Stochastic (online) truncated SVD for sparse matrices



Experiments: FOFE for Word Embedding

- Training Corpus: *enwiki9* (130 million words)
- Vocabulary: 80,000 words; Truncated dimension: 300
- FOFE matrices are weighted by **PMI**; Use **scipy** for truncated SVD
- Evaluated on five popular word similarity tasks:

Method	WordSim353	MEN	Mech Turk	Rare Words	SimLex-999
VSM+SVD	0.71	0.71	0.63	0.48	0.39
CBOW	0.68	0.68	0.66	0.43	0.35
SGNS	0.70	0.67	0.62	0.44	0.37
GloVe	0.59	0.64	0.58	0.39	0.29
Swivel	0.73	0.72	0.70	0.44	0.33
FOFE+SVD	0.76	0.76	0.65	0.50	0.39

Pearson correlation coefficients with human scores

FOFE-net for Language Modelling

 Use FOFE to *recursively* encode *partial sequence (history)* to predict next word



FOFE-net for Language Modelling

- Formulate FOFE as efficient matrix multiplications
- Much faster than RNN based language models

$$\mathbf{S} = \begin{bmatrix} 1 & & \\ \alpha & 1 & \\ \alpha^2 & \alpha & 1 & \\ \vdots & \ddots & 1 & \\ \alpha^{T-1} & \alpha & \alpha & 1 \end{bmatrix} \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \mathbf{e}_3 \\ \vdots \\ \mathbf{e}_T \end{bmatrix} = \mathbf{M}\mathbf{V}$$
$$\\ \mathbf{\bar{S}} = \begin{bmatrix} \mathbf{M}_1 & & \\ \mathbf{M}_2 & & \\ & \ddots & \\ & & \mathbf{M}_N \end{bmatrix} \begin{bmatrix} \mathbf{V}_1 \\ \mathbf{V}_2 \\ \vdots \\ \mathbf{V}_N \end{bmatrix} = \mathbf{\bar{M}}\mathbf{\bar{V}}.$$
$$\\ \mathbf{H} = f\left((\mathbf{\bar{M}}\mathbf{\bar{V}})\mathbf{U}\mathbf{W} + \mathbf{b}\right) = f\left(\mathbf{\bar{M}}(\mathbf{\bar{V}}\mathbf{U})\mathbf{W} + \mathbf{b}\right)$$

Experiments: FOFE for Language Modelling

- Penn Treebank (PTB) data set:
 - Training 930K words; valid 74K words; test 82K words
 - Vocabulary 10K words; *no drop-out*



Experiments: FOFE for Language Modelling

• Enwiki9 data set:

- Training 153M words; valid 8.9M words; test 8.9M words
- Vocabulary 80K words; *no drop-out*

Model	Architecture	Test PPL
KN 3-gram	-	156
KN 5-gram	-	132
	[1*200]-400-400-80k	241
	[2*200]-400-400-80k	155
FNN-LMs	[2*200]-600-600-80k	150
	[3*200]-400-400-80k	131
	[4*200]-400-400-80k	125
RNN-LM	[1*600]-80k	112
	[1*200]-400-400-80k	120
	[1*200]-600-600-80k	115
FORE FININ-LIVIS	[2*200]-400-400-80k	112
	[2*200]-600-600-80k	107



Local detection: no Viterbi decoding; Nested/Embedded entities

- No feature engineering: FOFE codes
- Easy and fast to train; make use of partial labels



Local detection: no Viterbi decoding; Nested/Embedded entities

- No feature engineering: FOFE codes
- Easy and fast to train; make use of partial labels



Local detection: no Viterbi decoding; Nested/Embedded entities

- No feature engineering: FOFE codes
- Easy and fast to train; make use of partial labels

- In training: labelled entities + negative sampling
- In test: consider all spans up to 7-8 words in a sentence



Experiments (I): FOFE-net for NER

Performance comparison on the CoNLL03 test set

FEATURE			Р	R	F1
		context FOFE incl. focus word(s)	86.64	77.04	81.56
	case-insensitive	context FOFE excl. focus word(s)	53.98	42.17	47.35
word level		BoW of focus word(s)	82.92	71.85	76.99
word-level		context FOFE incl. focus word(s)	88.88	79.83	84.12
	case-sensitive	context FOFE excl. focus word(s)	50.91	42.46	46.30
		BoW of focus word(s)	85.41	74.95	79.84
char FOFE of focus word(s)		cus word(s)	67.67	52.78	59.31
chai-level	Char CNN of focus word(s)			69.49	73.91
all case-inse	nsitive features		90.11	82.75	86.28
all case-sens	sitive features		90.26	86.63	88.41
all word-level features			92.03	86.08	88.96
all word-level & Char FOFE features			91.68	88.54	90.08
all word-level & Char CNN features			91.80	88.58	90.16
all word-level & all char-level features			93.29	88.27	90.71

	word	char	gaz	cap	pos	F1
(Collobert et al., 2011)	1	X	✓	✓	X	89.59
(Huang et al., 2015)	1	✓	✓	✓	✓	90.10
(Rondeau and Su, 2016)	1	X	✓	✓	✓	89.28
(Chiu and Nichols, 2016)	1	✓	✓	×	×	91.62
this work	1	✓	×	X	X	90.71

TAC-KBP EDL Contest

- Information extraction from unstructured text (Wikification)
- Use FOFE-net for entity discovery and mention detection



Experiments (II): 2015 KBP-EDL

- Pre-contest results: entity discovery performance comparison on 2015 KBP-EDL data set
- Trilingual: English, Chinese, Spanish
- 5 named entity types (PER, GPE, LOC, ORG, FAC) + 1 nominal mention (PER)

	2015 track best			ours			
	P	R	F_1	P	R	F_1	
Trilingual	75.9	69.3	72.4	78.3	69.9	73.9	
English	79.2	66.7	72.4	77.1	67.8	72.2	
Chinese	79.2	74.8	76.9	79.3	71.7	75.3	
Spanish	78.4	72.2	75.2	79.9	71.8	75.6	

Experiments (III): 2016 KBP-EDL official results

- 2016 KBP Trilingual EDL track: FOFE-net ranked No.2 among all participating teams
- 5 named entity types (PER, GPE, LOC, ORG, FAC) + all 5 nominal mentions

LANG	NAME		NOMINAL			OVERALL			
	Р	R	F1	Р	R	F1	Ρ	R	F1
	RUN1 (our official ED result in KBP2016 EDL2)								
ENG	0.898	0.789	0.840	0.554	0.336	0.418	0.836	0.680	0.750
CMN	0.848	0.702	0.768	0.414	0.258	0.318	0.789	0.625	0.698
SPA	0.835	0.778	0.806	0.000	0.000	0.000	0.835	0.602	0.700
ALL	0.893	0.759	0.821	0.541	0.315	0.398	0.819	0.639	0.718

FOFE-net for Coreference Resolution



Winograd Schema Challenge (WSC)

WSC: alternative Turing test for machine intelligence

The customer walked into the bank and stabbed one of the tellers. <u>**He</u>** was immediately taken to the **emergency**.</u>

Who was taken to the emergency room? The customer / the teller.

The customer walked into the bank and stabbed one of the tellers. <u>**He</u>** was immediately taken to the **police**.</u>

Who was taken to the emergency room? The customer / the teller.

2016 Winograd Schema Challenge

 Pronoun Disambiguation Problems (PDP): the first round qualifying test in 2016 WS challenge

Tom handed over the blueprints he had grabbed and, while his companion spread them out on <u>his</u> knee, walked toward the yard.

Tom / companion

One chilly May evening the English tutor invited Marjorie and myself into <u>her</u> room.

the English tutor / Marjorie

Mariano fell with a crash and lay stunned on the ground. Castello instantly kneeled by his side and raised <u>his</u> head.

Mariano / Castello

FOFE-net for PDP



Always before, **Larry** had helped Dad with his work. But he could not help him now, for Dad said that his boss at the railroad company would not want anyone but him to work in the office.



Incorporating Common-sense Knowledge



1. For each triple in ConceptNet:

 $(w_h, r, w_t) \Rightarrow sim(w_h, w_t) > sim(w_h, w_k) \quad w_k \in V \text{ and } w_k \text{ is not linked with } w_h$

2. For each cause-effect pair:

 $(w_i, w_j) \Rightarrow sim(w_i, w_j) > sim(w_i, w_k) \quad w_k \in V \text{ and } w_k \text{ is not the effect of } w_i$

Experiments (I): FOFE-net for PDP

 Official results of the top systems on the PDP test set of the 2016 Winograd Schema Challenge

Systems	Accuracy
Random guess	45%
Denis Robert	31.7%
Patrick Dhondt	45.0%
Nicos Issak	48.3%
Quan Liu	58.3 %

Experiments (II): FOFE-net for PDP

Post-contest results: more training data; more tuning

Taxt Corrous	Problem Solver	KEE settings Accuracy				
Text Corpus	FIODIeili Solvei	KEE training sources	Accuracy (%)	Improvements (%)		
		Context	48.3			
		Context + ConcepNet	55.0	+13.9		
	USSM	Context + WordNet	53.3	+10.4		
		Context + CauseCom	55.0	+13.9		
		Context + All KBs	56.7	+17.4		
		Context	<u>51.7</u>			
	NKAM	Context + ConcepNet	60.0	+16.0		
Wikipedia		Context + WordNet	60.0	+16.0		
		Context + CauseCom	61.7	+19.3		
		Context + All KBs	63.3	+22.4		
		Context	<u>53.3</u>			
		Context + ConcepNet	63.3	+18.7		
	USSM+NKAM	Context + WordNet	61.7	+15.7		
		Context + CauseCom	65.0	+21.9		
		Context + All KBs	66.7	+25.1		

Experiments (II): FOFE-net for PDP

Post-contest results: more training data; more tuning

Taxt Comput	Droblem Solver	KEE settings Accuracy				
Text Corpus	Problem Solver	KEE training sources	Accuracy (%)	Improvements (%)		
		Context	48.3			
		Context + ConcepNet	55.0	+13.9		
	USSM	Context + WordNet	53.3	+10.4		
		Context + CauseCom	55.0	+13.9		
		Context + All KBs	56.7	+17.4		
	NKAM	Context	<u>51.7</u>			
		Context + ConcepNet	60.0	+16.0		
Wikipedia		Context + WordNet	60.0	+16.0		
		Context + CauseCom	61.7	+19.3		
		Context + All KBs	63.3	+22.4		
		Context	<u>53.3</u>			
		Context + ConcepNet	63.3	+18.7		
	USSM+NKAM	Context + WordNet	61.7	+15.7		
		Context + CauseCom	65_0	+21.9		
		Context + All KBs	66.7	+25.1		

Outline

- Introduction
 - Deep Learning for NLP
- A New General Approach
 - FOFE: lossless encoding
 - DNNs: universal approximator
- NLP Applications
 - Word Embedding
 - Language Modelling
 - Named Entity Recognition and Mention Detection in KBP-EDL
 - Coreference Resolution in Winograd
- On-going work
 - Entity Linking
 - Syntactic/Semantic Parsing
- Conclusions



FOFE-net for Entity Linking

- Each named entity matches multiple items in knowledge bases
 - Name-Item matching score (matched char; FOFE-net)
 - Coherence between entities (max graph edges, personalized PageRank)



FOFE-net for Entity Linking



FOFE-net for Entity Linking





FOFE-net for Relation Extraction (Slot Filling)



FOFE-net for Syntactic Parsing

- Fit well with the transition-based parsing algorithm.
- No feature engineering; Fast to train from a very large corpus



Conclusions



- Proposed a new deep learning method (FOFE-net) for NLP
- FOFE-net:
 - Lossless fixed-size encoding + universal approximate
 - FOFE + neural network
 - Elegant and theoretically guaranteed
 - No feature engineering; Fast and easy to train
- FOFE-net is flexible: applicable to a wide range of NLP tasks
 - Word embedding
 - Language modelling
 - Named entity recognition
 - Coreference resolution
 - more and more ...
- Achieved promising results on all examined tasks