Natural Language Understanding in a Continuous Space

Karl-Moritz Hermann, Nal Kalchbrenner, Edward Grefenstette, and **Phil Blunsom**

phil.blunsom@cs.ox.ac.uk

Features and NLP



 $p(game|in) \propto exp(w^{T}\Phi(game,in))$ $\Phi_{1}(x,y) = \begin{cases} 1, & \text{if } PoS(x)=Noun \& y=in \\ 0, & \text{otherwise} \end{cases}$ $\Phi_{2}(x,y) = \begin{cases} 1, & \text{if } x=game \& PoS(y)=Prep \\ 0, & \text{otherwise} \end{cases}$ etc.

Twenty years ago log-linear models freed us from the shackles of simple multinomial parametrisations, but imposed the tyranny of feature engineering.



Distributed/neural models allow us to learn shallow features for our classifiers, capturing simple correlations between inputs.

Features and NLP



Deep learning allows us to learn hierarchical generalisations. Something that is proving rather useful for vision, speech, and now NLP...



2 From Vector Space Compositional Semantics to MT

We can represent words using a number of approaches

- Characters
- POS tags
- Grammatical roles
- Named Entity Recognition
- Collocation and distributional representations
- Task-specific features

All of these representations can be encoded in vectors. Some of these representations capture *meaning*.

A harder problem: paraphrase detection

Q: Do two sentences (roughly) mean the same? "He enjoys Jazz music" \equiv "He likes listening to Jazz" ?

A: Use a distributional representation to find out?

A harder problem: paraphrase detection

Q: Do two sentences (roughly) mean the same?"He enjoys Jazz music" ≡ "He likes listening to Jazz" ?

A: Use a distributional representation to find out?

Most representations not sensible on the sentence level

- Characters ?
- POS tags ?
- Grammatical roles ?
- Named Entity Recognition ?
- Collocation and distributional representations ?
- Task-specific features ?

The curse of dimensionality

As the dimensionality of a representation increases, learning becomes less and less viable due to sparsity.

Dimensionality for collocation

- One word per entry: Size of dictionary (small)
- One sentence per entry: Number of possible sentences (infinite)
- \Rightarrow We need a different method for representing sentences

Deep Learning for Language

Learning a hierarchy of features, where higher levels of abstraction are derived from lower levels.



A door, a roof, a window: It's a house



Composition

Lots of possible ways to compose vectors

- Addition
- Multiplication
- Kronecker Product
- Tensor Magic
- Matrix-Vector multiplication
- ...

Requirements

Not commutative Encode its parts? More than parts? Mary likes John \neq John likes Mary Magic carpet \equiv Magic + Carpet Memory lane \neq Memory + Lane

Autoencoders

We want to ensure that the joint representation captures the meaning of its parts. We can achieve this by autoencoding our data at each step:



For this to work, our autoencoder minimizes an objective function over inputs $x_i, i \in N$ and their reconstructions x'_i :

$$J = \frac{1}{2} \sum_{i}^{N} \left\| x_i' - x_i \right\|^2$$

We still want to learn how to represent a full sentence (or house). To do this, we chain autoencoders to create a recursive structure.



We use a composition function g(W * input + bias)

g is a non-linearity (tanh, sigm) W is a weight matrix b is a bias

A different task: paraphrase detection

Q: Do two sentences (roughly) mean the same? "He enjoys Jazz music" \equiv "He likes listening to Jazz" ?

A: Use deep learning to find out!



Other Applications: Stick a label on top



1. Combine label and reconstruction error

$$E(N, l, \theta) =$$

$$\sum_{n \in N} E_{rec}(n, \theta) + E_{lbl}(v_n, l, \theta)$$

$$E_{rec}(n, \theta) = \frac{1}{2} \left\| [x_n \| y_n] - r_n \right\|^2$$

$$E_{lbl}(v, l, \theta) = \frac{1}{2} \left\| l - v \right\|^2$$

2. Strong results for a number of tasks:

Sentiment Analysis Paraphrase Detection Image Search

. . .

Deep learning is suppose to learn the features for us, so can we do away with all this structural engineering and forget about latent parse trees?









A: My favourite show is Masterpiece Theatre.

A: Do you like it by any chance?

B: Oh yes!

A: You do!

B: Yes, very much.

A: Well, wouldn't you know.

B: As a matter of fact, I prefer public television.

B: And, uh, I have, particularly enjoy English comedies.

Statement-Non-Opinion Yes-No-Question Yes-Answers Declarative Yes-No-Q Yes-Answers Exclamation Statement-non-opinion Statement-non-opinion

Dave: Hello HAL, do you read me HAL?

HAL: Affirmative, Dave, I read you. Dave: Open the pod bay doors, HAL.

HAL: I'm sorry, Dave, I'm afraid I can't do that.



HAL: Affirmative, Dave, I read you.



Dave: Open the pod bay doors, HAL.



HAL: I'm sorry, Dave, I'm afraid I can't do that.







$$\begin{split} \mathbf{h}_i &= g(\mathbf{I}x_{i-1} + \mathbf{H}^{i-1}\mathbf{h}_{i-1} + \mathbf{Ss}_i)\\ p_i &= \mathsf{softmax}(\mathbf{O}^i\mathbf{h}_i) \end{split}$$



 $p_i = \mathsf{softmax}(\mathbf{O}^i \mathbf{h}_i)$

State of the art results while allowing online processing of dialogue.

Convolution Sentence Models: Question Answering



?x : have-population-of(vancouver, x)

Competitive with a template based approach with lots of hand engineered features.



The cat sat on the red mat













Small Sentiment Task

	Five-class (%)	Binary (%)
NB	41.0	81.8
SVM	40.7	79.4
BINB	41.9	83.1
RECNTN	45.7	85.4
MAX-TDNN	37.4	77.1
NBoW	42.4	80.5
DCNN	48.5	86.8

Sentiment prediction on the Stanford movie reviews dataset.

	Accuracy (%)
SVM	81.6
NB	82.7
MAXENT	83.0
MAX-TDNN	78.8
NBoW	80.9
DCNN	87.4

Accuracy on the larger Twitter sentiment dataset.

Question Classification Task

Classifier	Features	Acc. (%)
HIER	unigram, POS, head chunks NE, semantic relations	91.0
MAXENT	unigram, bigram, trigram POS, chunks, NE, supertags CCG parser, WordNet	92.6
MAXENT	unigram, bigram, trigram POS, wh-word, head word word shape, parser hypernyms, WordNet	93.6
SVM	unigram, POS, wh-word head word, parser hypernyms, WordNet 60 hand-coded rules	95.0
MAX-TDNN	unsupervised vectors	84.4
NBoW	unsupervised vectors	88.2
DCNN	unsupervised vectors	93.0

Six-way question classification on the TREC questions dataset, e.g.

Input: How far is it from Denver to Aspen ?

Output:

NUMBER

Feature: not only ... but also

not	only	manufactured	,	but	also	so
while	not	all	transitions	to	are	so
s	not	there	yet			but
not	all	transitions	to	are	so	,
may	not	be	new	,	but	australian
feels	not	only	manufactured	,	but	also
s	not	merely	unwatchable	,	but	also
land	than	crash	,	but	ultimately	serving
least	surprising	,	it	is	still	ultimately
great	bond	movie	,	but	it	is

as	predictable	and	as	lowbrow	as	the
as	lively	and	as	fun	as	it
,	confusing	spectacle	,	one	that	may
that	hinges	on	its	casting	,	and
cinematic	high	crime	,	one	that	brings
as	an	athlete	as	well	as	an
its	audience	and	its	source	material	
,	and	lane	as	vincent	,	the
as	lo	fi	as	the	special	effects
age	story	restraint	as	well	as	warmth

Feature: positivity

startling	film	that	gives	you	a	fascinating
well	written	,	nicely	acted	and	beautifully
best	from	his	large	cast	in	beautifully
strong	,	credible	performances	from	the	whole
compelling	journey		and	• •	his	best
be	a	joyful	or	at	least	fascinating
throughout	is	daring	,	inventive	and	impressive
enjoyable	film	for	the	family	,	amusing
originality	it	makes	up	for	in	intelligence
charming	,	quirky	and	paced	scottish	comedy



2 From Vector Space Compositional Semantics to MT

From Vector Space Compositional Semantics to MT

请给我一杯白葡萄酒。

From Vector Space Compositional Semantics to MT



i 'd like a glass of white wine , please .



i 'd like a glass of white wine , please .



Formal logical representations are very hard to learn from data. Let us optimistically assume a vector space and see how we go.

Generation

A simple distributed representation language model:



$$p_n = C_{n-2}R(w_{n-2}) + C_{n-1}R(w_{n-1})$$

$$p(w_n|w_{n-1}, w_{n-2}) \propto \exp(R(w_n)^T p_n)$$

This is referred to as a *log-bilinear model*.

Generation

A simple distributed representation language model:



$$p_n = C_{n-2}R(w_{n-2}) + C_{n-1}R(w_{n-1})$$
$$p(w_n|w_{n-1}, w_{n-2}) \propto \exp(R(w_n)^T \sigma(p_n))$$

Adding a non-linearity gives a version of what is often called a neural, or continuous space, LM.

Conditional Generation



$$p_n = C_{n-2}R(t_{n-2}) + C_{n-1}R(t_{n-1}) + \operatorname{CSM}(n, \mathbf{s})$$
$$p(t_n|t_{n-1}, t_{n-2}, \mathbf{s}) \propto \exp(R(t_n)^T \sigma(p_n))$$



$$p_n = C_2 R(t_{n-2}) + C_1 R(t_{n-1}) + \sum_{j=1}^{|\mathbf{s}|} S(s_j)$$
$$p(t_n | t_{n-1}, t_{n-2}, \mathbf{s}) \propto \exp(R(t_n)^T \sigma(p_n))$$

明天 早上 七点 叫醒 我 好 吗 ?

From Vector Space Compositional Semantics to MT







may i have a wake-up call at seven tomorrow morning ?







i 'm going to los angeles this afternoon .



i 'd like to have a room under thirty dollars a night .



i 'd like to have a room under thirty dollars a night .

Rough Gloss

I would like a night thirty dollars under room.



i 'd like to have a room under thirty dollars a night .

Google Translate

I want a late thirties under \$'s room.





Chinese (zh) \rightarrow English (en)	test 1	test 2
cdec (state-of-the-art MT)	50.1	58.9
Direct (naive bag of words source) Direct (convolution $p(en zh)$) Noisy Channel (convolution $p(zh en)p(en)$) Noisy Channel × Direct	30.8 44.6 50.1 51.0	33.2 50.4 51.8 55.2

BLEU score results on a small Chinese \rightarrow English translation task.

Advantages

- unsupervised features extraction alleviates domain and language dependencies.
- very compact models.
- distributed representations for words naturally include morphological properties.
- the conditional generation framework easily permits additional context such as dialogue and domain level vectors.

Challenges

- better conditioning on sentence position for long sentences, and all the other things this model does not capture!
- handling rare and unknown words.

Compositional Morphology for Word Representations and Language Modelling Botha and Blunsom. ICML 2014

Multilingual Models for Compositional Distributed Semantics Hermann and Blunsom. ACL 2014

A Convolutional Neural Network for Modelling Sentences Kalchbrenner, Grefenstette, and Blunsom. ACL 2014

Recurrent Continuous Translation Models Kalchbrenner and Blunsom. EMNLP 2013

The Role of Syntax in Vector Space Models of Compositional Semantics Hermann and Blunsom. ACL 2013



We are growing! Postdoctoral and DPhil studentships are available working in Machine Learning and Computational Linguistics

http://www.clg.ox.ac.uk

From Vector Space Compositional Semantics to MT