

Learning High-Order Word Representations

Konstantinos Kogkalidis

June 12, 2018

LoLa Fan Club

Motivation

Idea: structure-preserving map \mathcal{F}

$$\mathcal{F} : \mathcal{G} \rightarrow \mathbf{FdVect}$$

- Atomic types translated to vectors (order-one tensors)
- Complex types translated to (multi-)linear maps (higher order tensors)

Example

Word Type	\mathcal{G} Type	\mathcal{F} Translation
Noun	NP	\mathbb{R}^{NP}
Adjective	NP/NP	$\mathbb{R}^{NP \times NP} \equiv \mathbb{R}^{NP} \rightarrow \mathbb{R}^{NP}$

$$cat \in \mathbb{R}^{NP}$$

$$black, stray \in \mathbb{R}^{NP \times NP}$$

$$black\ stray\ cat \in \mathbb{R}^{NP}$$

Why Compositionality?

- Bridging of formal & distributional semantics
- Syntax-informed meaning derivations
- Modeling of functional words
- Formal treatment of ambiguous derivations
- Contextual Disambiguation
- Richer representations
- \vdots

Why Not Compositionality?

- ✓ Great properties
- ? How to obtain word representations?

Possible options:

Why Not Compositionality?

- ✓ Great properties
- ? How to obtain word representations?

Possible options:

1. Co-occurrence statistics

Why Not Compositionality?

✓ Great properties

? How to obtain word representations?

Possible options:

1. Co-occurrence statistics ✗
2. Unsupervised techniques (*a la word2vec*)

Why Not Compositionality?

- ✓ Great properties
- ? How to obtain word representations?

Possible options:

1. Co-occurrence statistics ✗
2. Unsupervised techniques (*a la word2vec*) ✗
3. Supervised learning ?

Problem Statement

Examine whether supervised learning can be used to find higher-order word representations (transitive verbs)

Supervised Learning

- Search over set of functions $A \rightarrow B$ parameterized over P
- Find optimal approximation \hat{f}_P to $f: A \rightarrow B$
- Use samples $(a, f(a)) \in A \times B$ to update P

Supervised Learning

Dataset

Sample space must be:

- labeled
- constrained
- of large size
- of high quality

Raw text paraphrase pairs

Example pair

proposed by the president ~ suggested by the chairman

- labeled ✓
- constrained ✗ (different syntactic types)
- of large size ✓
- of high quality ?

1. **Parse and filter by type** (transitive verb case)

- labeled ✓
- constrained ✓
- of large size ✗ (>95% loss)
- of high quality ✗ (*parser-induced errors*)

2. **Back-translation**

- Labeled ✓
- constrained ✓
- of large size ✓
- of high quality ✗ (*translation-induced errors*)

3. **Filter by co-occurrence / mutual information**

- Labeled ✓
- constrained ✓
- of large size ✓
- of high quality ?

Dataset: End Result

Verb / object dictionaries:

$$\mathcal{V} : \{v_1 : 1, v_2 : 2, \dots, v_N : N\}$$

$$\mathcal{O} : \{o_1 : 1, o_2 : 2, \dots, o_M : M\}$$

Paraphrase relation:

$$\mathcal{P} : \mathbb{N} \times \mathbb{N} \times \mathbb{N} \times \mathbb{N} \rightarrow \{0, 1\} \quad (\text{binary classification})$$

$$\mathcal{P}(i, j, k, l) = \mathcal{P}(k, l, i, j) = \begin{cases} 1 & v_i o_j \sim v_k o_l \\ 0 & \text{otherwise} \end{cases}$$

Supervised Learning

Formulating the Network

Training Objective

Our semantic interpretations are:

- Actions: $[a] = \mathbb{R}^A$
- Objects: $[np] = \mathbb{R}^{NP}$
- Transitive Verbs: $[a/np] = \mathbb{R}^{A \times NP}$

And our objective is to learn a **verb embedding function** ϵ_{verb} :

$$\epsilon_{verb} : \mathbb{N} \rightarrow \mathbb{R}^{A \times NP}$$

But instead we have samples from some $f : \mathbb{N}^4 \rightarrow \{0, 1\}$

Solution

Formulate f_p to incorporate ε_{verb} .

$$f_p = f_1 \circ f_2 \circ \dots \circ \varepsilon_{verb} \circ \dots$$

Intermediate Representations

Solution

Formulate f_p to incorporate ε_{verb} .

$$f_p = f_1 \circ f_2 \circ \dots \circ \varepsilon_{verb} \circ \dots$$

Simplification (1)

Assume pre-trained **object embedding function** ε_{object} .

$$\varepsilon_{object} : \mathbb{N} \rightarrow \mathbb{R}^{300}$$

Filling the missing blocks

$$i \in \mathbb{N}$$

$$j \in \mathbb{N}$$

$$k \in \mathbb{N}$$

$$l \in \mathbb{N}$$

$$\hat{y} \in \mathbb{R}$$

Filling the missing blocks

$$i \in \mathbb{N}$$

$$j \in \mathbb{N}$$

$$k \in \mathbb{N}$$

$$l \in \mathbb{N}$$

ε_{object}

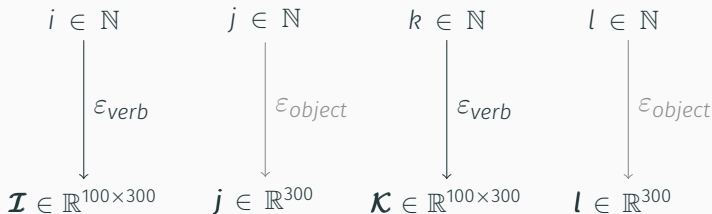
$$j \in \mathbb{R}^{300}$$

ε_{object}

$$l \in \mathbb{R}^{300}$$

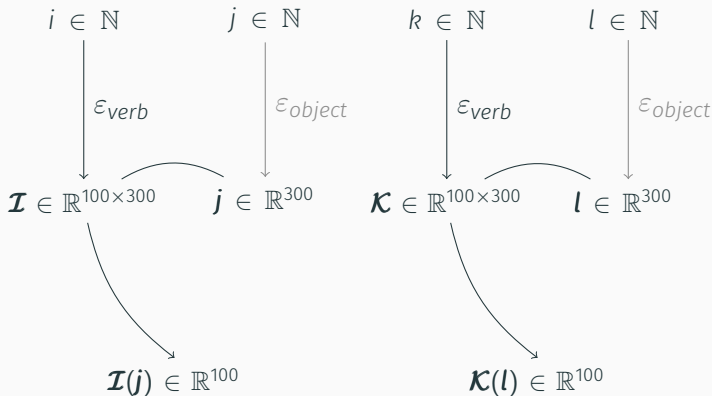
$$\hat{y} \in \mathbb{R}$$

Filling the missing blocks



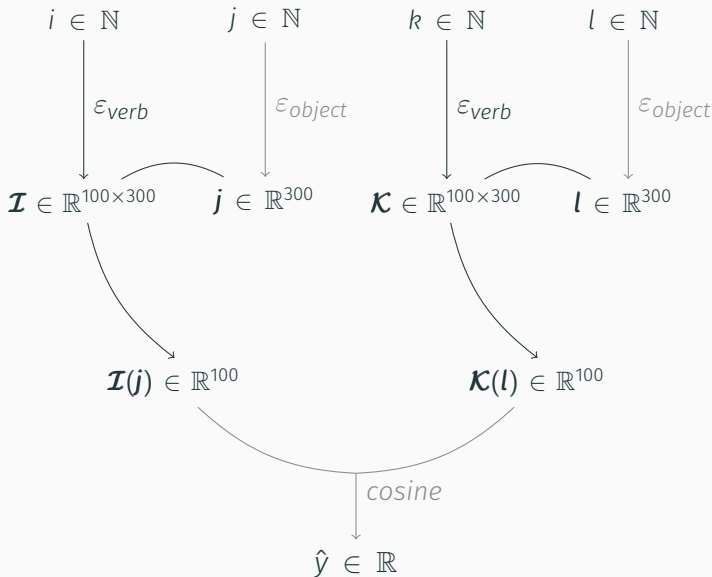
$$\hat{y} \in \mathbb{R}$$

Filling the missing blocks



$$\hat{y} \in \mathbb{R}$$

Filling the missing blocks



Objective Function

$$\cos(\mathbf{V}_i(j), \mathbf{V}_k(l)) \rightsquigarrow \mathcal{P}(i, j, k, l) \quad \forall (i, j, k, l) \in \mathcal{V} \times \mathcal{O} \times \mathcal{V} \times \mathcal{O}$$

Objective Function

$$\cos(\mathbf{V}_i(j), \mathbf{V}_k(l)) \rightsquigarrow \mathcal{P}(i, j, k, l) \quad \forall (i, j, k, l) \in \mathcal{V} \times \mathcal{O} \times \mathcal{V} \times \mathcal{O}$$

Considerations:

1. Network Size

- 1.000 verbs
 - $100 \times 300 = 30.000$ parameters per verb
- \Rightarrow 30.000.000 parameters for ε_{verb} to learn

Objective Function

$$\cos(\mathbf{V}_i(j), \mathbf{V}_k(l)) \rightsquigarrow \mathcal{P}(i, j, k, l) \quad \forall (i, j, k, l) \in \mathcal{V} \times \mathcal{O} \times \mathcal{V} \times \mathcal{O}$$

Considerations:

1. Network Size

- 1.000 verbs
 - $100 \times 300 = 30.000$ parameters per verb
- \Rightarrow 30.000.000 parameters for ε_{verb} to learn

2. Quantifying over two spaces ...

Objective Function

$$\cos(\mathbf{V}_i(j), \mathbf{V}_k(l)) \rightsquigarrow \mathcal{P}(i, j, k, l) \quad \forall (i, j, k, l) \in \mathcal{V} \times \mathcal{O} \times \mathcal{V} \times \mathcal{O}$$

Considerations:

1. Network Size

- 1.000 verbs
 - $100 \times 300 = 30.000$ parameters per verb
- \Rightarrow 30.000.000 parameters for ε_{verb} to learn

2. Quantifying over two spaces ...

3. ... both of which are non-convex



... A "beast" to train ☹️

Supervised Learning

Transferring Knowledge

Finding an Oracle

Dataception: use our labeled dataset to create a new dataset

Finding an Oracle

Dataception: use our labeled dataset to create a new dataset

Simplification (2)

Assume another pre-trained verb embedding function ϵ'_{verb} .

$$\epsilon'_{verb} : \mathbb{N} \rightarrow \mathbb{R}^{300}$$

Finding an Oracle

Dataception: use our labeled dataset to create a new dataset

Simplification (2)

Assume another pre-trained **verb embedding function** ϵ'_{verb} .

$$\epsilon'_{verb} : \mathbb{N} \rightarrow \mathbb{R}^{300}$$

Oracle

We can now train a **paraphrase embedding function** ϵ_{par} .

$$\epsilon_{par} : \mathbb{R}^{300} \times \mathbb{R}^{300} \rightarrow \mathbb{R}^{100}$$

$$i \in \mathbb{N}$$

$$j \in \mathbb{N}$$

$$k \in \mathbb{N}$$

$$l \in \mathbb{N}$$

$$\hat{y} \in \mathbb{R}$$

Oracle Flow

$$i \in \mathbb{N}$$

$$j \in \mathbb{N}$$

$$k \in \mathbb{N}$$

$$l \in \mathbb{N}$$

ε_{object}

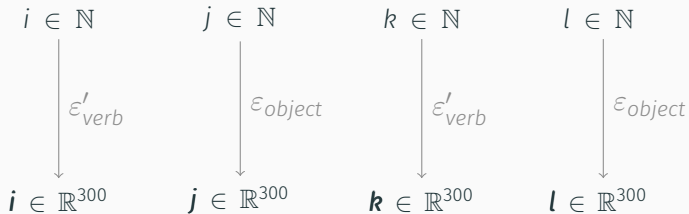
$$j \in \mathbb{R}^{300}$$

ε_{object}

$$l \in \mathbb{R}^{300}$$

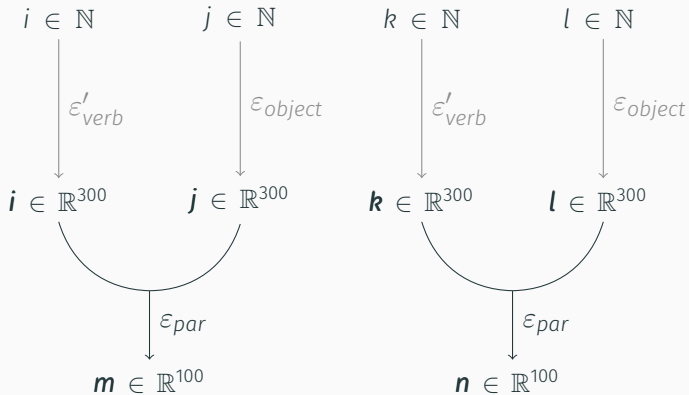
$$\hat{y} \in \mathbb{R}$$

Oracle Flow



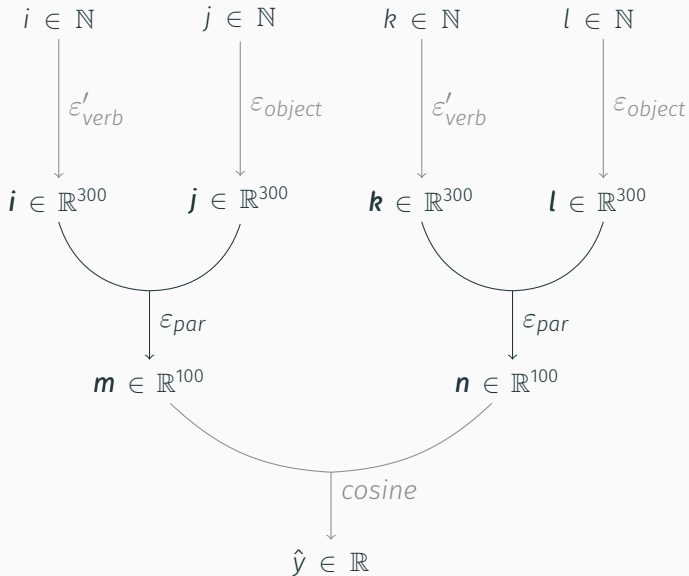
$$\hat{y} \in \mathbb{R}$$

Oracle Flow



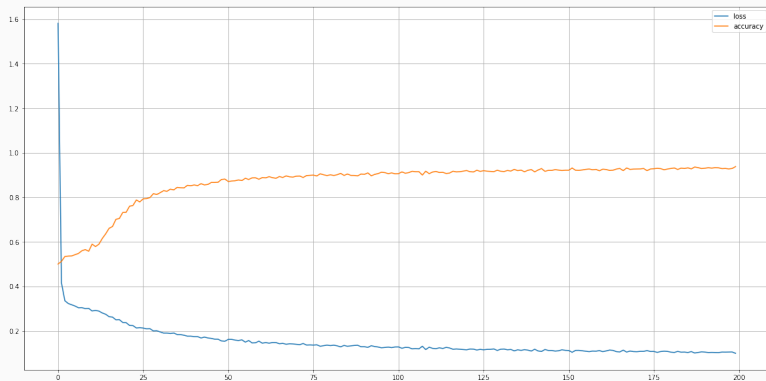
$$\hat{y} \in \mathbb{R}$$

Oracle Flow



Training the Oracle

ϵ_{par} : recurrent autoencoder (≈ 700.000 parameters)

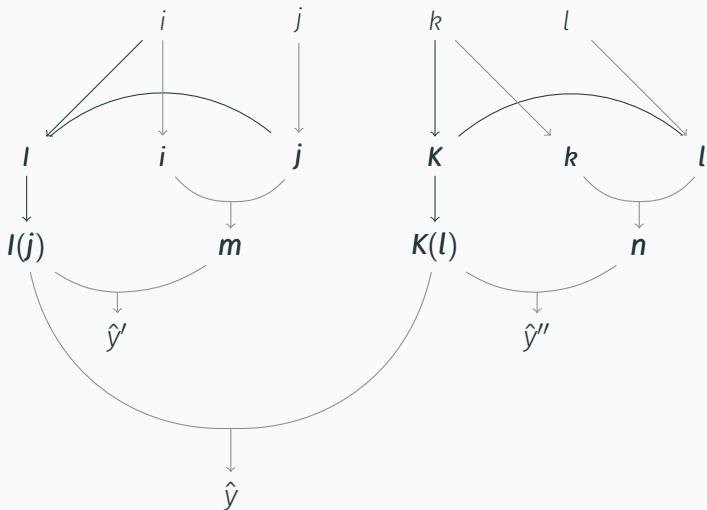


New Objective Function

$$\text{COS}(\mathbf{V}_i(\mathbf{j}), \varepsilon_{\text{par}}(\mathbf{v}_i, \mathbf{j})) \rightsquigarrow 1 \quad \forall (i, j) \in \mathcal{V} \times \mathcal{O}$$

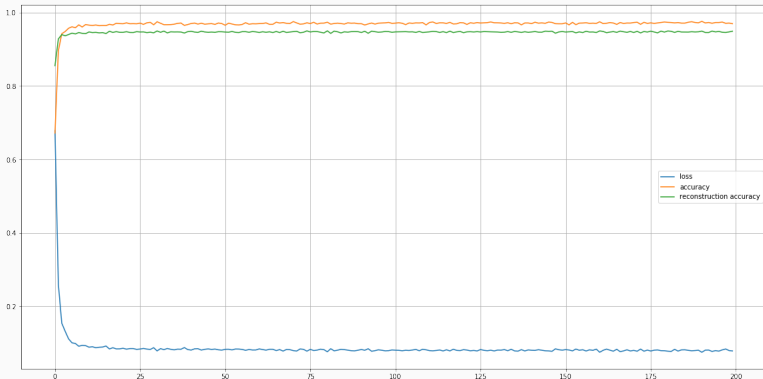
- ε_{par} gives us paraphrase embeddings '*for free*'
- We can use them to facilitate training
- Much smaller problem space

Composing Networks



Training the Original

ϵ_{verb} : tanh activated dense layer ($\approx 30.000.000$ parameters)





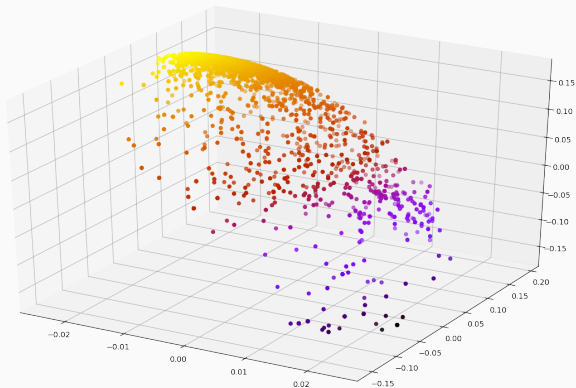
The beast has been tamed! 😊

Evaluation

Task-specific performance relates to the small-scale structure of the learned space:

Ground Truth \ Prediction	Oracle		Final	
	T	\perp	T	\perp
T	0.92	0.08	0.88	0.02
\perp	0.08	0.92	0.12	0.98

3D PCA on paraphrase embeddings



Conclusion

1. Metric reliability
2. Uninformative error signal
3. Over-parameterization
 - a) Linearity constraint
 - b) Chasing after an oracle
 - c) Bad scaling

Next Steps

1. Directly evaluate verb matrices
2. More structural constraints (activity regularization)
3. Iterative learning
4. Other data formats:
 - Different syntactic types
 - Different labels / samples altogether
5. Different oracle architectures
6. Different embedder architectures (encoder/decoder)

Bag of Tricks

Negative Sampling

Let $L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be the *loss function*.

Objective translates to:

$$\min_P L[\mathcal{P}(i, j, k, l), \hat{y}_P] \quad \forall (i, j, k, l) \quad (\text{Intractable})$$

Randomly generate and select negative samples (different every epoch).

"Two phrases are not similar unless they are"

- Class imbalance \Rightarrow bad predictions
- Treat negative samples as noise
- Learn positive examples then increase noise

Cross-entropy vs. MSE vs. Categorical Hinge

- Different assumptions, none correct.
- Many falsities \Rightarrow Truth ?