# Learning High-Order Word Representations

Konstantinos Kogkalidis

June 12, 2018

LoLa Fan Club

# Motivation

## Distributional Compositional Semantics

Idea: structure-preserving map $\mathcal{F}$

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

- 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 | $NP \backslash NP$ | $\mathbb{R}^{NP \times NP} \equiv \mathbb{R}^{NP} \to \mathbb{R}^{NP}$ |

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

## Why Compositionality?

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

## Why Not Compositionality?

✓ Great properties

**?** How to obtain word representations?

Possible options:

✓ Great properties

? How to obtain word representations?

Possible options:

1. Co-occurrence statistics

✓ Great properties

? How to obtain word representations?

Possible options:

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

✓ Great properties

**?** How to obtain word representations?

Possible options:

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

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 **?**

Verb / object dictionaries:

$$\mathcal{V} : \{v_1 : 1, v_2 : 2, ..., v_N : N\}$$
$$\mathcal{O} : \{o_1 : 1, o_2 : 2, ..., o_M : M\}$$

Paraphrase relation:

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

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

# Supervised Learning

## Formulating the Network

Our semantic interpretations are:

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

And our objective is to learn a verb embedding function $\varepsilon_{verb}$:

$$\varepsilon_{verb} : \mathbb{N} \to \mathbb{R}^{A \times NP}$$

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

**Solution**

Formulate $f_P$ to incorporate $\varepsilon_{verb}$.

$$f_P = f_1 \circ f_2 \circ \cdots \circ \varepsilon_{verb} \circ \ldots$$

## Solution

Formulate $f_P$ to incorporate $\varepsilon_{verb}$.

$$f_P = f_1 \circ f_2 \circ \cdots \circ \varepsilon_{verb} \circ \ldots$$

## Simplification (1)

Assume pre-trained object embedding function $\varepsilon_{object}$.

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

## 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

$$i \in \mathbb{N} \qquad j \in \mathbb{N} \qquad k \in \mathbb{N} \qquad l \in \mathbb{N}$$
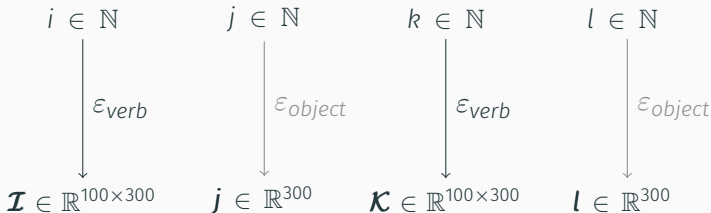
$\varepsilon_{object}$

$\varepsilon_{object}$

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

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

## Filling the missing blocks

$$i \in \mathbb{N} \qquad j \in \mathbb{N} \qquad k \in \mathbb{N} \qquad l \in \mathbb{N}$$

$$\downarrow \varepsilon_{verb} \qquad \downarrow \varepsilon_{object} \qquad \downarrow \varepsilon_{verb} \qquad \downarrow \varepsilon_{object}$$

$$\mathcal{I} \in \mathbb{R}^{100 \times 300} \quad j \in \mathbb{R}^{300} \quad \mathcal{K} \in \mathbb{R}^{100 \times 300} \quad l \in \mathbb{R}^{300}$$

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

# Filling the missing blocks

$$i \in \mathbb{N} \qquad j \in \mathbb{N} \qquad k \in \mathbb{N} \qquad l \in \mathbb{N}$$

$$\downarrow \varepsilon_{verb} \qquad \downarrow \varepsilon_{object} \qquad \downarrow \varepsilon_{verb} \qquad \downarrow \varepsilon_{object}$$

$$\mathcal{I} \in \mathbb{R}^{100 \times 300} \quad j \in \mathbb{R}^{300} \qquad \mathcal{K} \in \mathbb{R}^{100 \times 300} \quad l \in \mathbb{R}^{300}$$

$$\mathcal{I}(j) \in \mathbb{R}^{100} \qquad\qquad \mathcal{K}(l) \in \mathbb{R}^{100}$$

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

**Objective Function**

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

## Objective Function

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

Considerations:

1. Network Size
   - 1.000 verbs
   - $100 \times 300 = 30.000$ parameters per verb
   $\Rightarrow$ 30.000.000 parameters for $\varepsilon_{verb}$ to learn

## Objective Function

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

Considerations:

1. Network Size
   - 1.000 verbs
   - $100 \times 300 = 30.000$ parameters per verb
   $\Rightarrow$ 30.000.000 parameters for $\varepsilon_{verb}$ to learn

2. Quantifying over two spaces …

## Objective Function

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

Considerations:

1. Network Size
   - 1.000 verbs
   - $100 \times 300 = 30.000$ parameters per verb
   - $\Rightarrow$ 30.000.000 parameters for $\varepsilon_{verb}$ to learn

2. Quantifying over two spaces …

3. … both of which are non-convex

… A "beast" to train ☹

# Supervised Learning

## Transferring Knowledge

Dataception: use our labeled dataset to create a new dataset

Dataception: use our labeled dataset to create a new dataset

**Simplification (2)**

Assume another pre-trained verb embedding function $\varepsilon'_{verb}$.

$$\varepsilon'_{verb} : \mathbb{N} \to \mathbb{R}^{300}$$

Dataception: use our labeled dataset to create a new dataset

**Simplification (2)**

Assume another pre-trained verb embedding function $\varepsilon'_{verb}$.

$$\varepsilon'_{verb} : \mathbb{N} \to \mathbb{R}^{300}$$

**Oracle**

We can now train a paraphrase embedding function $\varepsilon_{par}$.

$$\varepsilon_{par} : \mathbb{R}^{300} \times \mathbb{R}^{300} \to \mathbb{R}^{100}$$

$i \in \mathbb{N}$ $j \in \mathbb{N}$ $k \in \mathbb{N}$ $l \in \mathbb{N}$

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

$$i \in \mathbb{N} \qquad j \in \mathbb{N} \qquad k \in \mathbb{N} \qquad l \in \mathbb{N}$$

$$\downarrow \varepsilon_{object} \qquad\qquad\qquad \downarrow \varepsilon_{object}$$

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

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

## Oracle Flow

$$i \in \mathbb{N} \qquad j \in \mathbb{N} \qquad k \in \mathbb{N} \qquad l \in \mathbb{N}$$

$$\downarrow \varepsilon'_{verb} \qquad \downarrow \varepsilon_{object} \qquad \downarrow \varepsilon'_{verb} \qquad \downarrow \varepsilon_{object}$$

$$\boldsymbol{i} \in \mathbb{R}^{300} \qquad \boldsymbol{j} \in \mathbb{R}^{300} \qquad \boldsymbol{k} \in \mathbb{R}^{300} \qquad \boldsymbol{l} \in \mathbb{R}^{300}$$

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

$i \in \mathbb{N}$    $j \in \mathbb{N}$    $k \in \mathbb{N}$    $l \in \mathbb{N}$

$\varepsilon'_{verb}$    $\varepsilon_{object}$    $\varepsilon'_{verb}$    $\varepsilon_{object}$

$\boldsymbol{i} \in \mathbb{R}^{300}$    $\boldsymbol{j} \in \mathbb{R}^{300}$    $\boldsymbol{k} \in \mathbb{R}^{300}$    $\boldsymbol{l} \in \mathbb{R}^{300}$

$\varepsilon_{par}$    $\varepsilon_{par}$

$\boldsymbol{m} \in \mathbb{R}^{100}$    $\boldsymbol{n} \in \mathbb{R}^{100}$

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

$i \in \mathbb{N}$  $j \in \mathbb{N}$  $k \in \mathbb{N}$  $l \in \mathbb{N}$

$\varepsilon'_{verb}$  $\varepsilon_{object}$  $\varepsilon'_{verb}$  $\varepsilon_{object}$

$\boldsymbol{i} \in \mathbb{R}^{300}$  $\boldsymbol{j} \in \mathbb{R}^{300}$  $\boldsymbol{k} \in \mathbb{R}^{300}$  $\boldsymbol{l} \in \mathbb{R}^{300}$

$\varepsilon_{par}$  $\varepsilon_{par}$

$\boldsymbol{m} \in \mathbb{R}^{100}$  $\boldsymbol{n} \in \mathbb{R}^{100}$

$cosine$

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

$\varepsilon_{par}$: recurrent autoencoder ($\approx 700.000$ parameters)
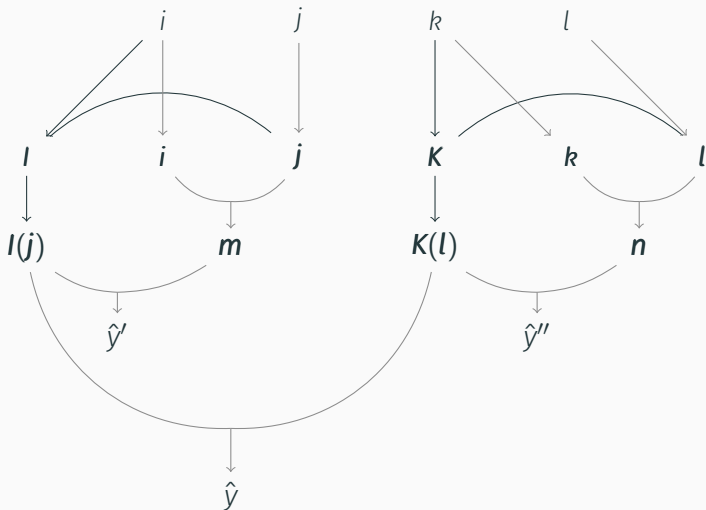
**New Objective Function**

$$cos(\mathbf{V}_i(\mathbf{j}), \varepsilon_{par}(\mathbf{v}, \mathbf{j})) \rightsquigarrow 1 \quad \forall \ (i, j) \in \mathcal{V} \times \mathcal{O}$$
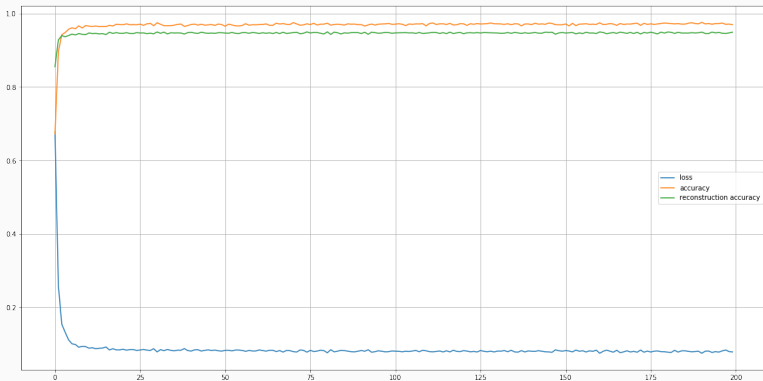
- $\varepsilon_{par}$ gives us paraphrase embeddings *'for free'*
- We can use them to facilitate training
- Much smaller problem space

$\varepsilon_{verb}$: tanh activated dense layer ($\approx$ 30.000.000 parameters)
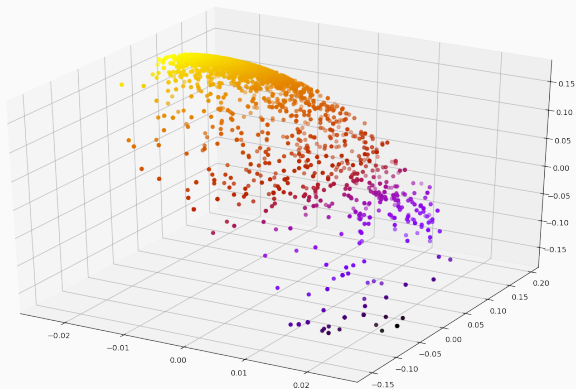
**The beast has been tamed!** ☺

# Evaluation

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

| Ground Truth \\ Prediction | Oracle $\top$ | Oracle $\bot$ | Final $\top$ | Final $\bot$ |
|---|---|---|---|---|
| $\top$ | 0.92 | 0.08 | 0.88 | 0.02 |
| $\bot$ | 0.08 | 0.92 | 0.12 | 0.98 |

## 3D PCA on paraphrase embeddings

# Conclusion

1. Uninformative error signal **?**
2. Over-parameterization
   a) Encoder / decoder architectures **?**
   b) Linearity constraint **✗**
   c) Chasing after an oracle **?**
   d) 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

# Bag of Tricks

1.  Negative Sampling
2.  Curriculum Learning / Domain Adaptation
3.  Loss Mixing