# Graph Embeddings for Natural Language Processing
## Lecture at Computer Science Club
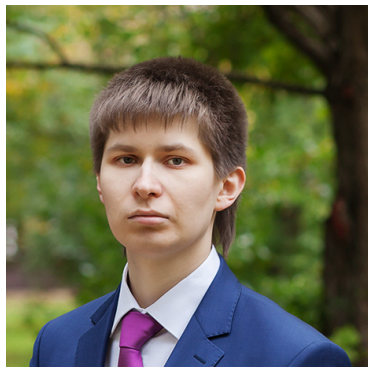
DOI: `10.5281/zenodo.4698904`

Dr. **Dmitry Ustalov**

Crowdsourcing Research Group
Yandex, Saint Petersburg, Russia

**Research Interests:** Crowdsourcing,
Computational Semantics, Evaluation

**Work Experience:** University of
Mannheim, Krasovskii Inst. of Math.
and Mech., Ural Federal University
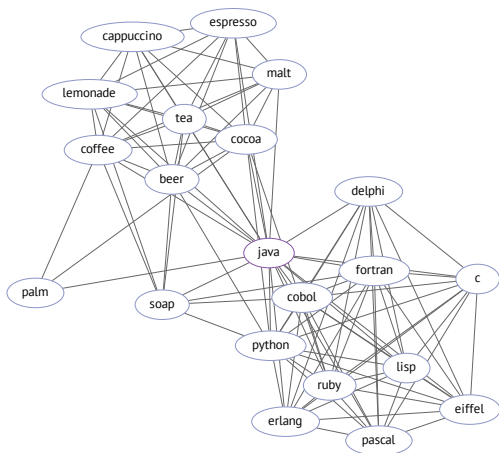
Section 1

Introduction

# Introduction

- Linguistic data are sparse, so the graphs are usually sparse, too
- Modern Natural Language Processing (NLP) is based on embeddings and representation learning
- We would like to reduce the dimensionality, but keep the important graph properties

$$\begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

### Core Idea: **Graphs are Features**

We can incorporate the relationships between objects in our machine learning pipelines.

# Motivation

Remember this *distributional thesaurus*?



- Can we measure the similarity between "tea" and "lisp"?
- Can we employ the node relationships as features?
- **Yes.**

Source: Ustalov et al. (2019)

# Successful Applications

Graph embeddings help in addressing very challenging NLP problems:

- question answering (Bordes et al., 2014)
- ranking for academic search (Xiong et al., 2017)
- text classification (Yao et al., 2019)
- fact checking (Zhong et al., 2020)
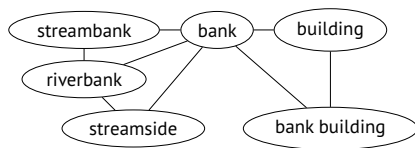- explanation regeneration (Li et al., 2020)

Beyond these applications, graph embeddings are generally useful for:

- node classification, recommendation, and link prediction
- feature extraction
- visualization (not every approach performs a proper layout)

# Problem Formulation

- There are node embeddings, edge embeddings, and the whole graph embeddings; we will focus on *node embeddings* only

- Given a graph $G = (V, E)$ and a number of dimensions $d \ll |V|$, we map $G$ into a $d$-dimensional space, in which the certain *graph property* is preserved as much as possible (Cai et al., 2018)

- Usually we would like to minimize some loss function using gradient descent

**Input Graph**



**Output Embedding**

# Section 2

## Unsupervised Embeddings

# Unsupervised Embeddings

- Unsupervised node embeddings build representations preserving generic graph properties
- We will focus on two different graph embedding methods: Laplacian Eigenmaps and DeepWalk
- There are *a lot* of other methods, see Cai et al. (2018) and Goyal et al. (2018)



Source: Finnsson (2017)

# Laplacian Eigenmaps (Spectral Embeddings)

- **Laplacian Eigenmaps** is a spectral approach for embedding high-dimensional data (Belkin et al., 2003)
- Compute a normalized Laplacian of the graph and run (approximate) *eigenvalue decomposition* to obtain the node embeddings
- Preserved graph properties are pairwise node similarities



Source: Amos (2011)

# Laplacian Eigenmaps: Algorithm

**Input:** graph $G = (V, E)$, adjacency matrix $A$, degree matrix $D$,
   dimensions $d \ll |V|$
**Output:** embedding $\vec{u} \in \mathbb{R}^d, \forall u \in V$
  1: $L^{\mathrm{norm}} \leftarrow D^{-\frac{1}{2}}(D - A)D^{-\frac{1}{2}}$
  2: $U\Lambda U^{-1} \leftarrow \mathtt{ein}(L^{\mathrm{norm}})$ $\quad$ ▷ Assume the eigenvalues are descending
  3: $U' \leftarrow (U_{ik})_{\substack{1 \leq i \leq |V|, 1 \leq j \leq d \\ k = |V|-1-j}}$ $\quad$ ▷ Drop the smallest eigenvalue
  4: **return** $\vec{u}_i \rightarrow U'_i$ **for all** $1 \leq i \leq |V|$

streambank

riverbank bank building

$$U' = \begin{pmatrix} .06 & 0 \\ -.31 & .71 \\ -.45 & 0 \\ -.31 & -.71 \\ .55 & 0 \\ .55 & 0 \end{pmatrix}$$

streamside

This is an example using the graph from Ustalov et al. (2019, Figure 2)

# Laplacian Eigenmaps: Discussion

Pros:
- ✚ Sound method that preserves local information optimally
- ✚ Very simple to implement

Cons:
- ➖ Slow, the worst-case running time is $O(|E|d^2)$
- ➖ Preserves only first-order proximity
- ➖ Graph should have only one connected component

Implementation:

🔗 https://scikit-learn.org/stable/modules/generated/sklearn.manifold.spectral_embedding.html

# Word2Vec Recap

- Mikolov et al. (2013) proposed Word2Vec, an efficient technique for learning *distributional representations* of words

- For each pair of word $w$ and its context $c$ in the fixed window, the Skip-Gram method performs negative sampling of $k \in \mathbb{N}$ contexts from $P_D$ and computes the objective (Levy et al., 2014):

$$\log \sigma(\vec{w} \cdot \vec{c}) + k \cdot \mathbb{E}_{c_N \sim P_D} \left[ \log \sigma(-\vec{w} \cdot \vec{c}_N) \right]$$

- Example representations:
$\overrightarrow{\text{Paris}} - \overrightarrow{\text{France}} + \overrightarrow{\text{Russia}} \approx \overrightarrow{\text{Moscow}}$
$\overrightarrow{\text{apple}} - \overrightarrow{\text{apples}} \approx \overrightarrow{\text{car}} - \overrightarrow{\text{cars}}$

- Popular variations are GloVe (Pennington et al., 2014), *fast*Text (Bojanowski et al., 2017), etc.

- ...but we are interested in graphs

# DeepWalk



- **DeepWalk** uses truncated random walks to learn latent representations by treating walks as the equivalent of *natural language sentences* (Perozzi et al., 2014)

- The input graph is flattened into a "corpus" of fixed-size node sequences; this corpus is used to train a Word2Vec model (Mikolov et al., 2013)
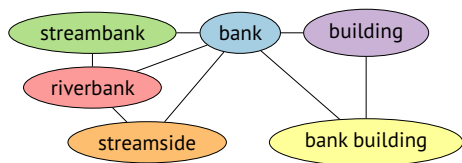
Source: Pexels (2016)

# DeepWalk: Algorithm

**Input:** graph $G = (V, E)$, dimensions $d \ll |V|$, window size $w \in \mathbb{N}$,
  walks per node $\gamma \in \mathbb{N}$, walk length $t \in \mathbb{N}$, learning rate $\alpha \in \mathbb{R}^+$
**Output:** embedding $\vec{u} \in \mathbb{R}^d, \forall u \in V$

1: $\Phi \leftarrow \texttt{random}(\mathbb{R}^{|V| \times d})$      ▷ Initialize from a uniform distribution
2: **for** $i \leftarrow 0 \ldots \gamma$ **do**
3:    **for all** $u \in V$ in random order **do**
4:      $\mathcal{W}_u \leftarrow \texttt{walk}(G, u, t)$      ▷ Random walk of length $t$ from $u$
5:      $\Phi \leftarrow \texttt{Skip-Gram}(\mathcal{W}_u, w, \alpha, \Phi)$      ▷ Update the parameters
6: **return** $\Phi$

 This is an example using the graph from Ustalov et al. (2019, Figure 2)

# DeepWalk: Discussion

## Pros:

**+** Very simple and works very well in practice

**+** Fast, the number of parameters is $O(d|V|)$

## Cons:

**−** Does not preserve community structure

**−** Does not preserve structural equivalence between nodes

**−** Edge weights are ignored

## Implementation:

🔗 https://github.com/phanein/deepwalk

🔗 https://snap.stanford.edu/node2vec/

🔗 http://rdf2vec.org/

## Word2Vec as Implicit Matrix Factorization

Levy et al. (2014) showed that Skip-Gram is an implicit factorization of a pointwise mutual information (PMI) word-context matrix.

- Given the word $w \in V$ and its context $c$, we count the number of words in context:

$$\text{PMI}(w, c) = \log \frac{\#(w, c) \cdot |D|}{\#(w) \cdot \#(c)}$$

- We obtain a *shifted PMI* matrix by shifting the PMI by a constant offset:

$$\text{SPPMI}_k(w, c) = \max(\text{PMI}(w, c) - \log k, 0)$$

- A truncated singular value decomposition $M^{\text{SPPMI}_k} = U_d \Sigma_d V_d^\top$ for the rank $d$ (Hansen, 1987) allows obtaining the embeddings $\Phi = U_d \sqrt{\Sigma_d}$ (here $V_d$ is a matrix and not a subset of $V$)

# "Embed All the Things!"

Wu et al. (2018) proposed a general-purpose embedding model **StarSpace**:

$$\sum_{(a,b) \in E^+} \sum_{b^- \in E^-} \underbrace{\max(0, \mu - \mathrm{sim}(a, b) + \mathrm{sim}(a, b^-))}_{\text{hinge loss with margin } \mu \in \mathbb{R}}$$

- Positive pairs $E^+$ are task-dependent and provided as the input; negative pairs $E^-$ are obtained by choosing $k \in \mathbb{N}$ negative pairs randomly
- Similarity function $\mathrm{sim}$ is either a dot product or cosine
- StarSpace is a convenient strong baseline for many tasks involving embedding entities comprised of discrete features: https://github.com/facebookresearch/StarSpace

# Unsupervised Embeddings: Wrap-Up

- Unsupervised node embeddings capture meaningful representations that can be concatenated or fine-tuned for downstream applications

- Edge weights can be accounted by performing graph traversal with BFS and DFS (Grover et al., 2016) or biased walks (Kartsaklis et al., 2018; Ristoski et al., 2018)



Source: rawpixel (2017)

# Section 3

## Supervised Embeddings

# Supervised Embeddings

- Building embeddings is not the ultimate goal: they are used in applications and there are useful features of the nodes
- **Graph Neural Networks** (GNNs) use the node features and relationships to learn node or graph representations
- We will focus on two semi-supervised graph embedding methods: Graph Convolutional Network and GraphSAGE
- There are *a lot* of others, see Dwivedi et al. (2020) and Wu et al. (2021)



Source: McGuire (2015)
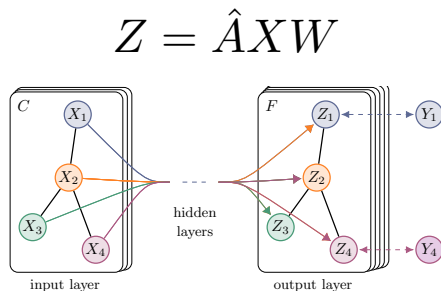
# Graph Neural Networks

- Given the graph $G = (V, E)$ and node feature vectors $X_v, \forall v \in V$, $k$-th layer of a GNN, $h_v^{(k)}$, is defined for node $v \in V$ (Xu et al., 2019):

$$
\left.
\begin{aligned}
a_v^{(k)} &= \text{aggregate}^{(k)} \left( \left\{ h_u^{(k-1)} : u \in V_v \right\} \right) \\
h_v^{(k)} &= \text{combine}^{(k)} \left( h_v^{(k-1)}, a_v^{(k)} \right)
\end{aligned}
\right\} \begin{array}{c} \text{differentiable} \\ \text{functions} \end{array}
$$

- After $k$ iterations, structural information about $k$-th order neighborhoods is captured; initialization is $h_v^{(0)} = X_v$

- Typically used with non-linear activation functions $\sigma(h_v^{(k)})$, such as $\tanh, \text{ReLU}, \text{softmax}$, etc.

❗ Parameters are estimated like in other kinds of neural networks, see Goodfellow et al. (2016)

# Graph Convolutional Network

- Kipf et al. (2017a) proposed a **Graph Convolutional Network** (GCN) that learns $F$-dimensional representations of graph nodes with $C$-dimensional features

- Given the signal $X \in \mathbb{R}^{|V| \times C}$, the convolved signal matrix is $Z = \hat{A} X W$, where $\hat{A}$ is the normalized adjacency matrix $A$ and $W \in \mathbb{R}^{C \times F}$ is the learned matrix of filter parameters

$$Z = \hat{A} X W$$



Source: Kipf et al. (2017b)

According to Xu et al. (2019):

$$h_v^{(k)} = \sigma \left( W \cdot \text{MEAN} \left\{ h_u^{(k-1)} : \forall u \in V_v \cup \{v\} \right\} \right)$$
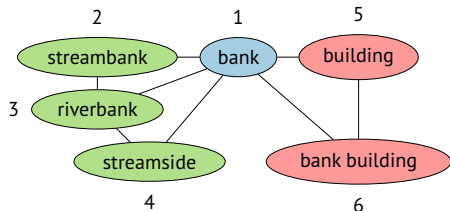
# Graph Convolutional Network: Estimation

As a semi-supervised method, GCN relies on labeled nodes $V_L \subseteq V$ and can be trained using the cross-entropy loss:

$$- \sum_{u_l \in V_L} \sum_{f=1}^{F} Y_{lf} \ln Z_{lf},$$

where $Y_{lf} = \begin{cases} 1, & \text{if } u_l \in V_L \text{ belongs to class } f, \\ 0, & \text{otherwise} \end{cases}$

- To avoid numerical instability, a *renormalization trick* with self-loops is used: $\tilde{A} = A + I$ and $\tilde{D}_{ii} = \sum_{1 \leq j \leq |V|} \tilde{A}_{ij}$, so $\hat{A} = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$
- Note that $F \ll |V|$ is the target number of dimensions $d$

This is an example using the graph from Ustalov et al. (2019, Figure 2)

# Graph Convolutional Network: Discussion

Pros:

➕ Sound method that approximates localized spectral filters on graphs

➕ Fast, the running time is linear in the number of edges

Cons:

➖ Prone to over-smoothing (Chen et al., 2020)

➖ Exact algorithm requires the complete $\hat{A}$

Implementations:

🔗 https://github.com/tkipf/gcn

🔗 https://github.com/tkipf/pygcn

# GraphSAGE

Hamilton et al. (2017) proposed a method for *inductive* node embedding that performs sampling and feature aggregation (GraphSAGE).

According to Xu et al. (2019):

$$a_v^{(k)} = \text{MAX} \left( \left\{ \sigma \left( W \cdot h_u^{(k-1)} \right) : \forall u \in V_v \right\} \right)$$

$$h_v^{(k)} = \sigma \left( W \cdot \left( h_v^{(k-1)} \oplus a_v^{(k)} \right) \right)$$
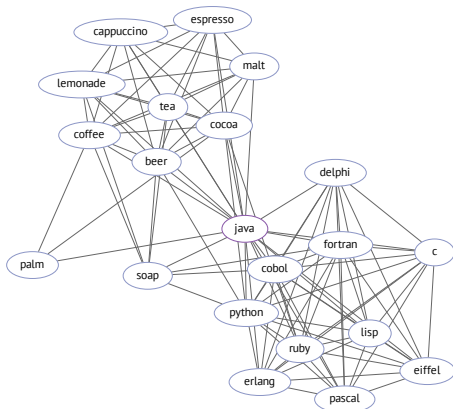
- Due to the sampling it allows not keeping the whole adjacency matrix in memory
- GraphSAGE generalizes over the GNN approach to use trainable aggregation functions and provides an unsupervised setting

# GraphSAGE: Estimation

GraphSAGE offers an unsupervised loss function for every representation $\vec{v} \in \mathbb{R}^d, \forall v \in V$ using $Q \in \mathbb{N}$ nodes sampled from the negative sampling distribution $P_n(v)$:

$$-\underbrace{\log\left(\sigma(\vec{u}^\top \vec{v})\right)}_{\text{adjacent nodes}} - \underbrace{Q \cdot \mathbb{E}_{v_n \sim P_n(v)} \log\left(\sigma(-\vec{u}^\top \vec{v_n})\right)}_{Q \text{ negative samples}}$$

- Representations of nearby nodes are similar,
  while representations of distant nodes are different
- Can be augmented with a task-specific loss,
  e.g., cross-entropy like in Kipf et al. (2017a)

# GraphSAGE: Example



| Node 1 | Node 2 | Cosine |
|--------|--------|--------|
| java | tea | .91 |
| java | lisp | −.35 |
| ruby | lisp | .99 |
| tea | coffee | .95 |
| tea | lisp | −.68 |

| Node | Vector |
|------|--------|
| java | $(.09, -.12, .22, -.02)^\top$ |
| tea | $(.09, -.15, .18, .08)^\top$ |
| lisp | $(-.11, .35, -.03, -.37)^\top$ |
| ruby | $(-.06, .29, .01, -.30)^\top$ |

This is an example using the graph from Ustalov et al. (2019, Figure 2)

# GraphSAGE: Discussion

Pros:

+ Sound method that approximates clustering coefficients
+ Allows differentiable aggregators, e.g., LSTMs (Hochreiter et al., 1997)
+ Allows unsupervised training and inductive setting

Cons:

− Still prone to over-smoothing (Chen et al., 2020)
− More hyper-parameters for tuning

Implementations:

🔗 https://github.com/williamleif/GraphSAGE
🔗 https://github.com/williamleif/graphsage-simple

# Supervised Embeddings: Wrap-Up

- Node embeddings can be efficiently estimated by the specific task
- These representations can be learned and extracted from the neural networks
- Semi-supervised representations do not require the complete data annotation
- Even a single layer of a GNN improves quality in practice (we will look at case studies)



Source: FreePhotosART (2016)

Section 4

Case Studies

# Case Studies

- Embedding a Distributional Thesaurus (Jana et al., 2018)
- Mapping Text to Knowledge Graphs (Kartsaklis et al., 2018)
- Semantic Role Labeling (Marcheggiani et al., 2017)
- Explanation Regeneration (Jansen et al., 2020)

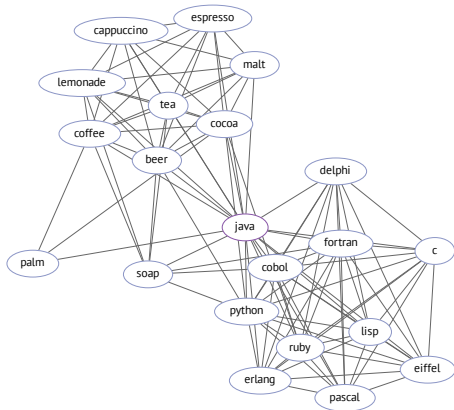

Source: Simone_ph (2017)

# Embedding DTs

- Jana et al. (2018) used embeddings of nodes in a distributional thesaurus (DT) as additional features for building better word representations



Source: Buissinne (2016)

# Embedding DTs: Approach

1. Build a distributional thesaurus (Biemann et al., 2013)
2. Learn node embeddings (DeepWalk, node2vec, etc.)
3. Concatenate node embeddings with GloVe word embeddings (Pennington et al., 2014)
4. Perform a principal component analysis (PCA)



Source: Ustalov et al. (2019)
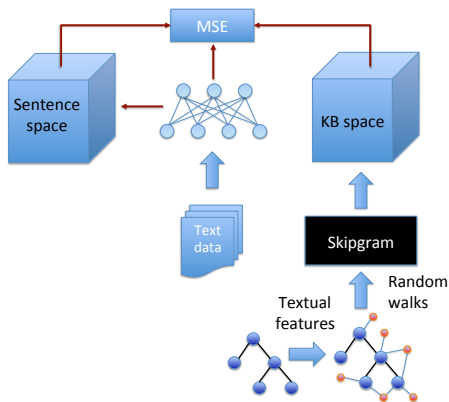
# Embedding DTs: Results

- According to the Spearman's $\rho$, concatenation (CC) of GloVe vectors with the DeepWalk embeddings improved the results on multiple datasets

- Note that PCA also improved upon CC despite the loss of information while dimensionsality reduction from 300 + 128 to 300

| Dataset | GloVe | CC | PCA |
|---------|-------|-------|-------|
| WSSim | 0.799 | 0.838 | 0.839 |
| SimL-N | 0.427 | 0.443 | 0.468 |
| RG-65 | 0.791 | 0.816 | **0.879** |
| MC-30 | 0.799 | 0.860 | **0.890** |
| WSR | 0.637 | **0.676** | 0.645 |
| M771 | 0.707 | 0.708 | 0.707 |
| M287 | 0.800 | 0.781 | 0.807 |
| MEN-N | **0.819** | 0.792 | 0.799 |
| WS-353 | 0.706 | **0.751** | 0.740 |

Source: Jana et al. (2018)
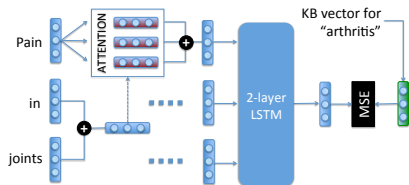
# Text-to-Entity Mapping

- Kartsaklis et al. (2018) proposed a technique for enriching the entity vectors with textual information
- Textual information is obtained from BabelNet (Navigli et al., 2012) and other sources



Source: Kartsaklis et al. (2018)

# Text-to-Entity Mapping: Approach

1. Learn node embeddings with DeepWalk (Perozzi et al., 2014)

2. Build LSTM (Hochreiter et al., 1997) with multi-sense aspect (aka MS-LSTM)

3. Minimize the mean squared error (MSE) between the sense vector and the target entity vector



Source: Kartsaklis et al. (2018)

Code and Data: https://bitbucket.org/dimkart/ms-lstm

# Text-to-Entity Mapping: Example

table[1]     formulation, uncommonly, rauwolfia, cardiology, hypodermic, malleability, points, optic, dendrite, rubiaceae, nonparametric, meninges, deviation, anesthetics

table[2]     tableware, meal, expectation, heartily, kitchen, hum, eating, forestay, suitors, croupier, companionship, restaurant, dishes, candles, cup, tea

table[3]     reassigned, projective, ultracentrifuge, polemoniaceous, thyronine, assumptions, lymphocyte, atomic, difficulties, intracellular, virgil, elementary, cartesian
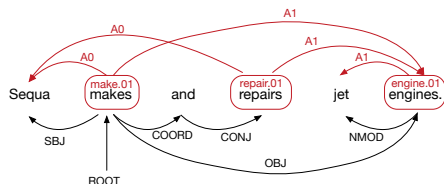
Source: Kartsaklis et al. (2018)

# Text-to-Entity Mapping: Results

- On the SMOMED CT dataset the text-to-entity mapping outperforms Word2Vec-based baselines
- On reverse dictionary and node classification tasks it shows results comparable to the state-of-the-art techniques (Kartsaklis et al., 2018)

| Model | Target | Accuracy |
|---|---|---|
| Baseline | W2V-GoogleNews | 0.19 |
| | W2V-PubMed | 0.12 |
| MS-LSTM | DeepWalk | 0.26 |
| | Enhanced | **0.84** |

Source: Kartsaklis et al. (2018)

- **Semantic Role Labeling** (SRL) assigns to the words in sentence the labels corresponding to their semantic role
- Marcheggiani et al. (2017) is the first paper that demonstrates the effectiveness of GCNs for NLP in the SRL setup
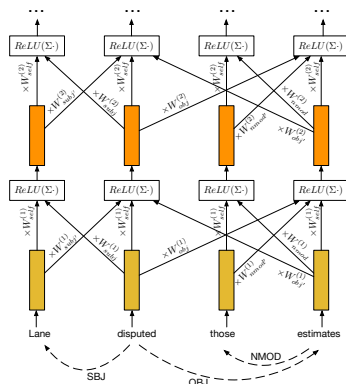


Source: Marcheggiani et al. (2017)

- Syntactic dependency trees are directed, so the layer is
$$h_v^{(k+1)} = \sigma \left( \sum_{u \in V_v} g_{v,u}^{(k)} (V_{\text{dir}(u,v)}^{(k)} h_u^{(k)} + b_{L(u,v)}^{(k)}) \right)$$

- For each edge-node pair there is a scalar gate:
$$g_{u,v}^{(k)} = \sigma \left( h_u^{(k)} \cdot \hat{v}_{\text{dir}(u,v)}^{(k)} + \hat{b}_{L(u,v)}^{(k)} \right)$$
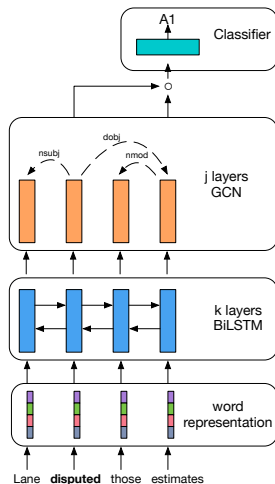


Source: Marcheggiani et al. (2017)

# GCNs for SRL: Approach



1. Fetch word embeddings
2. Stack several BiLSTM layers (Hochreiter et al., 1997)
3. Stack several GCN layers (Kipf et al., 2017a)
4. Add a softmax classifier

Code and Data: https://github.com/diegma/neural-dep-srl
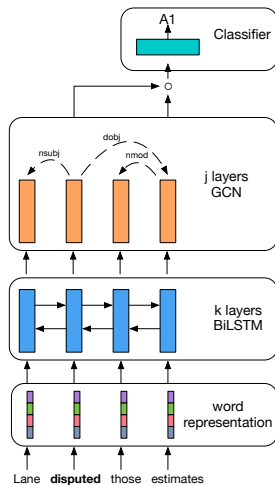
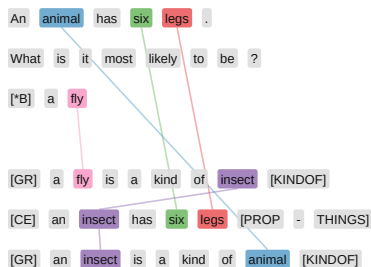Source: Marcheggiani et al. (2017)

# GCNs for SRL: Results

- GCN for SRL outperformed other approaches on both English and Chinese on the CoNLL-2009 dataset
- LSTMs without GCNs outperform GCNs without LSTMs, while their combination dramatically improves the precision
- Even a single GCN layer increases the LSTM-based model accuracy



Source: Marcheggiani et al. (2017)

# Explanation Regeneration

- In the **Explanation Regeneration** task, given an elementary science question with an answer to it, one has to rank explanations of this answer (Jansen et al., 2020)

- The best-performing system at the TextGraphs-14 shared task combined language models and graph netural networks (Li et al., 2020)



Source: Jansen et al. (2020)

# Explanation Regeneration: Approach

1. Retrieve the relevant explanations for the questions using ERNIE 2.0 (Sun et al., 2020)

2. Re-rank the retrieved sentences using ERNIE 2.0

3. Aggregate them using the GraphSAGE-like approach (Hamilton et al., 2017)



Source: Li et al. (2020)

Code and Data: https://github.com/PaddlePaddle/PGL/tree/static_stable/examples/erniesage

# Explanation Regeneration: Example

**?** A student placed an ice cube on a plate in the sun.
Ten minutes later, only water was on the plate.
Which process caused the ice cube to change to water?
(A) condensation     (B) evaporation     (C) freezing     (D) melting

| Rank | Gold | Fact (Table Row) |
|---|---|---|
| 1 | ⋆ | melting is a kind of process |
| 2 | | thawing is similar to melting |
| 3 | | melting is a kind of phase change |
| 4 | | melting is when solids are heated above their melting point |
| 5 | | amount of water in a body of water increases by (storms ; rain ; ice melting) |
| 6 | | an ice cube is a kind of object |
| 7 | ⋆ | an ice cube is a kind of solid |
| 8 | | freezing point is similar to melting point |
| 9 | | melting point is a property of a (substance ; material) |
| 10 | | glaciers melting has a negative impact on the glaicial environment |

...

Source: Jansen et al. (2020)

# Explanation Regeneration: Results

- According to Mean Average Precision (MAP), all the systems have dramatically improved over the tf-idf baseline
- Other systems used BERT, LSTM, integer linear programming, but the best system, BPGL, combined *texts and graphs* (Li et al., 2020)

| Model | MAP |
|-------|------|
| tf-idf | 0.23 |
| AG | 0.37 |
| RDAI | 0.55 |
| CSX | 0.50 |
| LIIR | 0.57 |
| BPGL | **0.60** |

Source: Jansen et al. (2020)

Section 5

Conclusion

# Conclusion

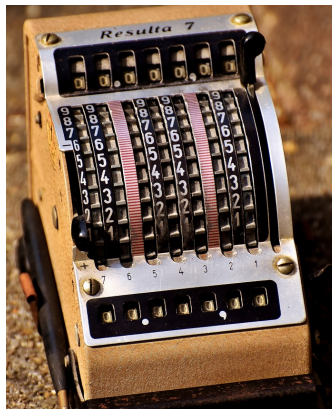- Node embeddings allow incorporating relationships between nodes in machine learning pipeline
- These techniques improve quality and are available in unsupervised, semi-supervised, and fully supervised setups
- Not covered here: knowledge graph embeddings (Wang et al., 2017), interpretability (Şenel et al., 2018), relationships with BERT-like models (Devlin et al., 2019)



Source: bamenny (2016)

# Implementations



- PyTorch Geometric (PyG)
  (Fey et al., 2019)
- PGL (Ma et al., 2019)
- DGL (Wang et al., 2019)
- GraphGym (You et al., 2020)
- Karate Club (Rozemberczki et al., 2020)

Source: Alexas_Fotos (2017)

# Questions?

### Contacts

## Dr. **Dmitry Ustalov**

Crowdsourcing Research Group
Yandex, Saint Petersburg, Russia

 https://github.com/dustalov

 mailto:dmitry.ustalov@gmail.com

 0000-0002-9979-2188

Revision: 47c51af

# References I

Belkin M. and Niyogi P. (2003). Laplacian Eigenmaps for Dimensionality Reduction and Data Representation. *Neural Computation*, vol. 15, no. 6, pp. 1373–1396. DOI: 10.1162/089976603321780317.

Biemann C. and Riedl M. (2013). Text: now in 2D! A framework for lexical expansion with contextual similarity. *Journal of Language Modelling*, vol. 1, no. 1, pp. 55–95. DOI: 10.15398/jlm.v1i1.60.

Bojanowski P. et al. (2017). Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 135–146. DOI: 10.1162/tacl_a_00051.

Bordes A., Chopra S., and Weston J. (2014). Question Answering with Subgraph Embeddings. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. EMNLP 2014. Doha, Qatar: Association for Computational Linguistics, pp. 615–620. DOI: 10.3115/v1/D14-1067.

Cai H., Zheng V. W., and Chen-Chuan Chang K. (2018). A Comprehensive Survey of Graph Embedding: Problems, Techniques, and Applications. *IEEE Transactions on Knowledge and Data Engineering*, vol. 30, no. 9, pp. 1616–1637. DOI: 10.1109/TKDE.2018.2807452.

Chen D. et al. (2020). Measuring and Relieving the Over-Smoothing Problem for Graph Neural Networks from the Topological View. *Proceedings of the AAAI Conference on Artificial Intelligence*, AAAI-20 vol. 34, no. 5, pp. 3438–3445. DOI: 10.1609/aaai.v34i04.5747.

Devlin J. et al. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. NAACL-HLT 2019. Minneapolis, MN, USA: Association for Computational Linguistics, pp. 4171–4186. DOI: 10.18653/v1/N19-1423.

Dwivedi V. P. et al. (2020). Benchmarking Graph Neural Networks. arXiv: 2003.00982 [cs.LG].

Fey M. and Lenssen J. E. (2019). Fast Graph Representation Learning with PyTorch Geometric. *ICLR 2019 Workshop on Representation Learning on Graphs and Manifolds*. arXiv: 1903.02428 [cs.LG]. URL: https://rlgm.github.io/papers/2.pdf.

Goodfellow I., Bengio Y., and Courville A. (2016). Deep Learning. Cambridge, MA, USA: MIT Press. ISBN: 978-0-262-03561-3. URL: https://www.deeplearningbook.org/.

Goyal P. and Ferrara E. (2018). Graph embedding techniques, applications, and performance: A survey. *Knowledge-Based Systems*, vol. 151, pp. 78–94. DOI: 10.1016/j.knosys.2018.03.022.

Grover A. and Leskovec J. (2016). node2vec: Scalable Feature Learning for Networks. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '16. San Francisco, CA, USA: ACM, pp. 855–864. DOI: 10.1145/2939672.2939754.

# References II

Hamilton W. L., Ying R., and Leskovec J. (2017). Inductive Representation Learning on Large Graphs. *Advances in Neural Information Processing Systems 30*. NIPS 2017. Vancouver, BC, Canada: Curran Associates, Inc., pp. 1024–1034. URL: `https://proceedings.neurips.cc/paper/2017/file/5dd9db5e033da9c6fb5ba83c7a7ebea9-Paper.pdf`.

Hansen P. C. (1987). The truncated *SVD* as a method for regularization. *BIT Numerical Mathematics*, vol. 27, no. 4, pp. 534–553. DOI: `10.1007/BF01937276`.

Hochreiter S. and Schmidhuber J. (1997). Long Short-Term Memory. *Neural Computation*, vol. 9, no. 8, pp. 1735–1780. DOI: `10.1162/neco.1997.9.8.1735`.

Jana A. and Goyal P. (2018). Can Network Embedding of Distributional Thesaurus Be Combined with Word Vectors for Better Representation? *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. NAACL-HLT 2018. New Orleans, LA, USA: Association for Computational Linguistics, pp. 463–473. DOI: `10.18653/v1/N18-1043`.

Jansen P. and Ustalov D. (2020). TextGraphs 2020 Shared Task on Multi-Hop Inference for Explanation Regeneration. *Proceedings of the Graph-based Methods for Natural Language Processing (TextGraphs)*. Barcelona, Spain (Online): Association for Computational Linguistics, pp. 85–97. URL: `https://www.aclweb.org/anthology/2020.textgraphs-1.10`.

Kartsaklis D., Pilehvar M. T., and Collier N. (2018). Mapping Text to Knowledge Graph Entities using Multi-Sense LSTMs. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. EMNLP 2018. Brussels, Belgium: Association for Computational Linguistics, pp. 1959–1970. DOI: `10.18653/v1/D18-1221`.

Kipf T. N. and Welling M. (2017a). Semi-Supervised Classification with Graph Convolutional Networks. *5th International Conference on Learning Representations, Conference Track Proceedings*. ICLR 2017. Toulon, France: OpenReview.net. URL: `https://openreview.net/forum?id=SJU4ayYgl`.

Levy O. and Goldberg Y. (2014). Neural Word Embedding as Implicit Matrix Factorization. *Advances in Neural Information Processing Systems 27*. NIPS 2014. Montréal, QC, Canada: Curran Associates, Inc., pp. 2177–2185. URL: `https://proceedings.neurips.cc/paper/2014/file/feab05aa91085b7a8012516bc3533958-Paper.pdf`.

Li W. et al. (2020). PGL at TextGraphs 2020 Shared Task: Explanation Regeneration using Language and Graph Learning Methods. *Proceedings of the Graph-based Methods for Natural Language Processing (TextGraphs)*. Barcelona, Spain (Online): Association for Computational Linguistics, pp. 98–102. URL: `https://www.aclweb.org/anthology/2020.textgraphs-1.11`.

Ma Y. et al. (2019). PaddlePaddle: An Open-Source Deep Learning Platform from Industrial Practice. *Frontiers of Data and Computing*, vol. 1, no. 1, pp. 105–115. DOI: `10.11871/jfdc.issn.2096.742X.2019.01.011`.

# References III

Marcheggiani D. and Titov I. (2017). Encoding Sentences with Graph Convolutional Networks for Semantic Role Labeling. *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. EMNLP 2017. Copenhagen, Denmark: Association for Computational Linguistics, pp. 1506–1515. DOI: `10.18653/v1/D17-1159`.

Mikolov T. et al. (2013). Distributed Representations of Words and Phrases and their Compositionality. *Advances in Neural Information Processing Systems 26*. NIPS 2013. Lake Tahoe, NV, USA: Curran Associates, Inc., pp. 3111–3119. URL: `https://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf`.

Navigli R. and Ponzetto S. P. (2012). BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artificial Intelligence*, vol. 193, pp. 217–250. DOI: `10.1016/j.artint.2012.07.001`.

Pennington J., Socher R., and Manning C. (2014). GloVe: Global Vectors for Word Representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. EMNLP 2014. Doha, Qatar: Association for Computational Linguistics, pp. 1532–1543. DOI: `10.3115/v1/D14-1162`.

Perozzi B., Al-Rfou R., and Skiena S. (2014). DeepWalk: Online Learning of Social Representations. *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '14. New York, NY, USA: ACM, pp. 701–710. DOI: `10.1145/2623330.2623732`.

Ristoski P. et al. (2018). RDF2Vec: RDF graph embeddings and their applications. *Semantic Web*, pp. 1–32. DOI: `10.3233/SW-180317`.

Rozemberczki B., Kiss O., and Sarkar R. (2020). Karate Club: An API Oriented Open-Source Python Framework for Unsupervised Learning on Graphs. *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. CIKM '20. Virtual Event, Ireland: Association for Computing Machinery, pp. 3125–3132. DOI: `10.1145/3340531.3412757`.

Şenel L. K. et al. (2018). Semantic Structure and Interpretability of Word Embeddings. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 10, pp. 1769–1779. DOI: `10.1109/TASLP.2018.2837384`.

Sun Y. et al. (2020). ERNIE 2.0: A Continual Pre-Training Framework for Language Understanding. *Proceedings of the AAAI Conference on Artificial Intelligence*, AAAI-20 vol. 34, no. 5, pp. 8968–8975. DOI: `10.1609/aaai.v34i05.6428`.

Ustalov D. et al. (2019). Watset: Local-Global Graph Clustering with Applications in Sense and Frame Induction. *Computational Linguistics*, vol. 45, no. 3, pp. 423–479. DOI: `10.1162/COLI_a_00354`.

Wang M. et al. (2019). Deep Graph Library: Towards Efficient And Scalable Deep Learning on Graphs. *ICLR 2019 Workshop on Representation Learning on Graphs and Manifolds*. URL: `https://rlgm.github.io/papers/49.pdf`.

Wang Q. et al. (2017). Knowledge Graph Embedding: A Survey of Approaches and Applications. *IEEE Transactions on Knowledge and Data Engineering*, vol. 29, no. 12, pp. 2724–2743. DOI: `10.1109/TKDE.2017.2754499`.

# References IV

Wu L. et al. (2018). StarSpace: Embed All The Things! *The Thirty-Second AAAI Conference on Artificial Intelligence*. AAAI-18. Association for the Advancement of Artificial Intelligence, pp. 5569–5577. URL: https://ojs.aaai.org/index.php/AAAI/article/view/11996.

Wu Z. et al. (2021). A Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 1, pp. 4–24. DOI: 10.1109/TNNLS.2020.2978386.

Xiong C., Power R., and Callan J. (2017). Explicit Semantic Ranking for Academic Search via Knowledge Graph Embedding. *Proceedings of the 26th International Conference on World Wide Web*. WWW '17. Perth, WA, Australia: International World Wide Web Conferences Steering Committee, pp. 1271–1279. DOI: 10.1145/3038912.3052558.

Xu K. et al. (2019). How Powerful are Graph Neural Networks? *7th International Conference on Learning Representations, Conference Track Proceedings*. ICLR 2019. New Orleans, LA, USA: OpenReview.net. URL: https://openreview.net/forum?id=ryGs6iA5Km.

Yao L., Mao C., and Luo Y. (2019). Graph Convolutional Networks for Text Classification. *Proceedings of the AAAI Conference on Artificial Intelligence*, AAAI-19, IAAI-19, EAAI-20 vol. 33, no. 1, pp. 7370–7377. DOI: 10.1609/aaai.v33i01.33017370.

You J., Ying Z., and Leskovec J. (2020). Design Space for Graph Neural Networks. *Advances in Neural Information Processing Systems 33*. NeurIPS 2020. Montréal, QC, Canada: Curran Associates, Inc., pp. 17009–17021. URL: https://proceedings.neurips.cc/paper/2020/file/c5c3d4fe6b2cc463c7d7ecba17cc9de7-Paper.pdf.

Zhong W. et al. (2020). Reasoning Over Semantic-Level Graph for Fact Checking. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. ACL 2020. Online: Association for Computational Linguistics, pp. 6170–6180. DOI: 10.18653/v1/2020.acl-main.549.

# Supplementary Media I

Alexas_Fotos (October 7, 2017). Calculating Machine Resulta Old. Pixabay. URL: https://pixabay.com/images/id-2825179/. Licensed under Pixabay License.

Amos E. (December 19, 2011). The Vectrex video game console, shown with controller. Wikimedia Commons. URL: https://commons.wikimedia.org/wiki/File:Vectrex-Console-Set.jpg. Licensed under CC BY-SA 3.0, used with author's permission.

bamenny (February 24, 2016). Robot Flower Technology. Pixabay. URL: https://pixabay.com/images/id-1214536/. Licensed under Pixabay License.

Buissinne S. (August 25, 2016). Dictionary Reference Book Learning. Pixabay. URL: https://pixabay.com/images/id-1619740/. Licensed under Pixabay License.

Finnsson I. (May 19, 2017). Books Covers Book Case. Pixabay. URL: https://pixabay.com/images/id-2321934/. Licensed under Pixabay License.

FreePhotosART (September 3, 2016). Cook Cooking School Pan. Pixabay. URL: https://pixabay.com/images/id-1641959/. Licensed under Pixabay License.

Kipf T. N. and Welling M. (February 22, 2017b). Semi-Supervised Classification with Graph Convolutional Networks. arXiv: 1609.02907v4 [cs.LG]. Licensed under arXiv.org perpetual, non-exclusive license, used with author's permission.

McGuire R. (March 24, 2015). Suit Business Man. Pixabay. URL: https://pixabay.com/images/id-673697/. Licensed under Pixabay License.

Pexels (November 23, 2016). Aquarium Jellyfish Aquatic. Pixabay. URL: https://pixabay.com/images/id-1851643/. Licensed under Pixabay License.

rawpixel (April 18, 2017). Calm Freedom Location. Pixabay. URL: https://pixabay.com/images/id-2218409/. Licensed under Pixabay License.

Simone_ph (March 21, 2017). Music Low Electric Bass. Pixabay. URL: https://pixabay.com/images/id-2149880/. Licensed under Pixabay License.