Deep Learning on Graphs with Graph Convolutional Networks



Thomas Kipf, 6 April 2017

joint work with Max Welling (University of Amsterdam)



UNIVERSITEIT VAN AMSTERDAM



The success story of deep learning

IM GENET







The success story of deep learning

IM AGENET





Deep neural nets that exploit:

- translation invariance (weight sharing)
- hierarchical compositionality





Deep Learning on Graph-Structured Data

Euclidean data: grids, sequences...







Euclidean data: grids, sequences...



Deep Learning on Graph-Structured Data

Convolutional neural networks (CNNs)





Convolutional neural networks (CNNs)





Convolutional neural networks (CNNs)



Recurrent neural networks (RNNs)



Traditional vs. "deep" learning

Traditional approach





Traditional vs. "deep" learning

Traditional approach



End-to-end learning



Single CNN layer with 3x3 filter:





Single CNN layer with 3x3 filter:





Single CNN layer with 3x3 filter:







Single CNN layer with 3x3 filter:





Update for a single pixel:

- Transform messages individually $\mathbf{W}_i \mathbf{h}_i$
- Add everything up $\sum_i \mathbf{W}_i \mathbf{h}_i$

Single CNN layer with 3x3 filter:





Update for a single pixel:

- Transform messages individually $\mathbf{W}_i \mathbf{h}_i$
- Add everything up $\sum_i \mathbf{W}_i \mathbf{h}_i$

Full update:
$$\mathbf{h}_{4}^{(l+1)} = \sigma \left(\mathbf{W}_{0}^{(l)} \mathbf{h}_{0}^{(l)} + \mathbf{W}_{1}^{(l)} \mathbf{h}_{1}^{(l)} + \dots + \mathbf{W}_{8}^{(l)} \mathbf{h}_{8}^{(l)} \right)$$

What if our data looks like this?













What if our data looks like this?





Real-world examples:

- Social networks
- World-wide-web
- Protein-interaction networks
- Telecommunication networks
- Knowledge graphs

• ...

Graph: $G = (\mathcal{V}, \mathcal{E})$

Adjacency matrix: A





Graph: $G = (\mathcal{V}, \mathcal{E})$

Adjacency matrix: A





Model wish list:

- Trainable in $\mathcal{O}(|\mathcal{E}|)$ time
- Applicable even if the input graph changes

- Take adjacency matrix ${f A}$ and feature matrix ${f X}$ •
- Concatenate them $\mathbf{X}_{in} = [\mathbf{A}, \mathbf{X}]$ •
- Feed them into deep (fully connected) neural net •



9

- Take adjacency matrix ${\bf A}$ and feature matrix ${\bf X}$
- Concatenate them $\mathbf{X}_{\mathrm{in}} = [\mathbf{A}, \mathbf{X}]$
- Feed them into deep (fully connected) neural net



Problems:

- Huge number of parameters $\mathcal{O}(N)$
- Re-train if graph changes

- Take adjacency matrix ${\bf A}$ and feature matrix ${\bf X}$
- Concatenate them $\mathbf{X}_{\mathrm{in}} = [\mathbf{A}, \mathbf{X}]$
- Feed them into deep (fully connected) neural net



Problems:

- Huge number of parameters $\mathcal{O}(N)$
- Re-train if graph changes

- Take adjacency matrix ${\bf A}$ and feature matrix ${\bf X}$
- Concatenate them $\mathbf{X}_{\mathrm{in}} = [\mathbf{A}, \mathbf{X}]$
- Feed them into deep (fully connected) neural net



Problems:

- Huge number of parameters $\mathcal{O}(N)$
- Re-train if graph changes

(related idea was first proposed in Scarselli et al. 2009)

Consider this undirected graph:





(related idea was first proposed in Scarselli et al. 2009)

Consider this undirected graph:



Calculate update for node in red:



(related idea was first proposed in Scarselli et al. 2009)

Consider this undirected graph:



Calculate update for node in red:



(related idea was first proposed in Scarselli et al. 2009)

Consider this
undirected graph:Calculate update
for node in red:Image: Consider this
for node in red:Image: Constant of the second sec

Update
rule:
$$\mathbf{h}_{i}^{(l+1)} = \sigma \left(\mathbf{h}_{i}^{(l)} \mathbf{W}_{0}^{(l)} + \sum_{j \in \mathcal{N}_{i}} \frac{1}{c_{ij}} \mathbf{h}_{j}^{(l)} \mathbf{W}_{1}^{(l)} \right) \begin{array}{l} \mathcal{N}_{i} : \text{neighbor indices} \\ c_{ij} : \text{norm. constant} \\ (\text{per edge}) \end{array}$$

Note: We could also choose simpler or more general functions over the neighborhood

Deep Learning on Graph-Structured Data

GCN model architecture

Input: Feature matrix $\mathbf{X} \in \mathbb{R}^{N imes E}$, preprocessed adjacency matrix $\hat{\mathbf{A}}$



[Kipf & Welling, ICLR 2017]

What does it do? An example.

Forward pass through untrained 3-layer GCN model



A "classical" approach for node feature assignment



Algorithm 1: WL-1 algorithm (Weisfeiler & Lehmann, 1968)

Input: Initial node coloring $(h_1^{(0)}, h_2^{(0)}, ..., h_N^{(0)})$ Output: Final node coloring $(h_1^{(T)}, h_2^{(T)}, ..., h_N^{(T)})$ t $\leftarrow 0$; repeat for $v_i \in \mathcal{V}$ do $\left[\begin{array}{c} h_i^{(t+1)} \leftarrow hash\left(\sum_{j \in \mathcal{N}_i} h_j^{(t)}\right); \\ t \leftarrow t+1; \end{array} \right]$ until stable node coloring is reached;

13

A "classical" approach for node feature assignment



Algorithm 1: WL-1 algorithm (Weisfeiler & Lehmann, 1968)

Input: Initial node coloring $(h_1^{(0)}, h_2^{(0)}, ..., h_N^{(0)})$ Output: Final node coloring $(h_1^{(T)}, h_2^{(T)}, ..., h_N^{(T)})$ t $\leftarrow 0$; repeat for $v_i \in \mathcal{V}$ do $\left[\begin{array}{c} h_i^{(t+1)} \leftarrow hash\left(\sum_{j \in \mathcal{N}_i} h_j^{(t)}\right); \\ t \leftarrow t+1; \end{array} \right]$ until stable node coloring is reached;

A "classical" approach for node feature assignment



Algorithm 1: WL-1 algorithm (Weisfeiler & Lehmann, 1968) Input: Initial node coloring $(h_1^{(0)}, h_2^{(0)}, ..., h_N^{(0)})$

Input: Initial node coloring $(h_1^{(T)}, h_2^{(T)}, ..., h_N^{(T)})$ Output: Final node coloring $(h_1^{(T)}, h_2^{(T)}, ..., h_N^{(T)})$ $t \leftarrow 0;$ repeat for $v_i \in \mathcal{V}$ do $\left[\begin{array}{c} h_i^{(t+1)} \leftarrow hash\left(\sum_{j \in \mathcal{N}_i} h_j^{(t)}\right); \\ t \leftarrow t+1; \end{array} \right]$ until stable node coloring is reached;

Useful as graph isomorphism check for most graphs

(exception: highly regular graphs)

A "classical" approach for node feature assignment



Algorithm 1: WL-1 algorithm (Weisfeiler & Lehmann, 1968) Input: Initial node coloring $(h_1^{(0)}, h_2^{(0)}, ..., h_N^{(0)})$ Output: Final node coloring $(h_1^{(T)}, h_2^{(T)}, ..., h_N^{(T)})$ t $\leftarrow 0$; repeat for $v_i \in \mathcal{V}$ do $\left[\begin{array}{c} h_i^{(t+1)} \leftarrow \operatorname{bach}\left(\sum_{j \in \mathcal{N}_i} h_j^{(t)}\right); \\ t \leftarrow t+1; \end{array} \right]$ until stable node coloring is reached;

Useful as graph isomorphism check for most graphs

(exception: highly regular graphs)

A "classical" approach for node feature assignment



Useful as graph isomorphism check for most graphs

(exception: highly regular graphs)

Setting:

Some nodes are labeled (black circle) All other nodes are unlabeled

Task:

Predict node label of unlabeled nodes





Setting:

Some nodes are labeled (black circle) All other nodes are unlabeled

Task:

Predict node label of unlabeled nodes



Standard approach:

graph-based regularization [Zhu et al., 2003]

$$\mathcal{L} = \mathcal{L}_0 + \lambda \mathcal{L}_{\text{reg}}$$
 with $\mathcal{L}_{\text{reg}} = \sum_{i,j} A_{ij} \|f(X_i) - f(X_j)\|^2$

assumes: connected nodes likely to share same label

Embedding-based approaches

Two-step pipeline:

- 1) Get embedding for every node
- 2) Train classifier on node embedding

Examples: DeepWalk [Perozzi et al., 2014], node2vec [Grover & Leskovec, 2016]



Embedding-based approaches

Two-step pipeline:

- 1) Get embedding for every node
- 2) Train classifier on node embedding

Examples: DeepWalk [Perozzi et al., 2014], node2vec [Grover & Leskovec, 2016]

Problem: Embeddings are not optimized for classification!



Embedding-based approaches

Two-step pipeline:

- 1) Get embedding for every node
- 2) Train classifier on node embedding

Examples: DeepWalk [Perozzi et al., 2014], node2vec [Grover & Leskovec, 2016]

Problem: Embeddings are not optimized for classification!

Idea: Train graph-based classifier end-to-end using GCN

Evaluate loss on labeled nodes only:

$$\mathcal{L} = -\sum_{l \in \mathcal{Y}_L} \sum_{f=1}^F Y_{lf} \ln Z_{lf}$$

- \mathcal{Y}_L set of labeled node indices
- ${f Y}$ label matrix
- ${f Z}\,$ GCN output (after softmax)

Toy example (semi-supervised learning)



Deep Learning on Graph-Structured Data

Toy example (semi-supervised learning)



Deep Learning on Graph-Structured Data

Application: Classification on citation networks

Input: Citation networks (nodes are papers, edges are citation links, optionally bag-of-words features on nodes)

Target: Paper category (e.g. stat.ML, cs.LG, ...)



Experiments and results

Model: 2-layer GCN $Z = f(X, A) = \operatorname{softmax}\left(\hat{A} \operatorname{ReLU}\left(\hat{A}XW^{(0)}\right)W^{(1)}\right)$

Dataset statistics

| Dataset | Туре | Nodes | Edges | Classes | Features | Label rate |
|----------|------------------|--------|---------|---------|----------|------------|
| Citeseer | Citation network | 3,327 | 4,732 | 6 | 3,703 | 0.036 |
| Cora | Citation network | 2,708 | 5,429 | 7 | 1,433 | 0.052 |
| Pubmed | Citation network | 19,717 | 44,338 | 3 | 500 | 0.003 |
| NELL | Knowledge graph | 65,755 | 266,144 | 210 | 5,414 | 0.001 |

(Kipf & Welling, Semi-Supervised Classification with Graph Convolutional Networks, ICLR 2017)

Experiments and results

Model: 2-layer GCN $Z = f(X, A) = \operatorname{softmax}\left(\hat{A} \operatorname{ReLU}\left(\hat{A}XW^{(0)}\right)W^{(1)}\right)$

Dataset statistics

| Dataset | Туре | Nodes | Edges | Classes | Features | Label rate |
|----------|------------------|--------|---------|---------|----------|------------|
| Citeseer | Citation network | 3,327 | 4,732 | 6 | 3,703 | 0.036 |
| Cora | Citation network | 2,708 | 5,429 | 7 | 1,433 | 0.052 |
| Pubmed | Citation network | 19,717 | 44,338 | 3 | 500 | 0.003 |
| NELL | Knowledge graph | 65,755 | 266,144 | 210 | 5,414 | 0.001 |

Classification results (accuracy)

| | Method | Citeseer | Cora | Pubmed | NELL |
|-------------------|--------------------|------------------|--------------|-------------------|----------------|
| | ManiReg [3] | 60.1 | 59.5 | 70.7 | 21.8 |
| | SemiEmb [24] | 59.6 | 59.0 | 71.1 | 26.7 |
| | → LP [27] | 45.3 | 68.0 | 63.0 | 26.5 |
| no input features | DeepWalk [18] | 43.2 | 67.2 | 65.3 | 58.1 |
| | Planetoid* [25] | 64.7 (26s) | 75.7 (13s) | 77.2 (25s) | 61.9 (185s) |
| | GCN (this paper) | 70.3 (7s) | 81.5 (4s) | 79.0 (38s) | 66.0 (48s) |
| | GCN (rand. splits) | 67.9 ± 0.5 | 80.1 ± 0.5 | 78.9 ± 0.7 | 58.4 ± 1.7 |

(Kipf & Welling, Semi-Supervised Classification with Graph Convolutional Networks, ICLR 2017)

Experiments and results

Model: 2-layer GCN $Z = f(X, A) = \operatorname{softmax}\left(\hat{A} \operatorname{ReLU}\left(\hat{A}XW^{(0)}\right)W^{(1)}\right)$



(Kipf & Welling, Semi-Supervised Classification with Graph Convolutional Networks, ICLR 2017)

Other recent applications

Molecules



[Duvenaud et al., NIPS 2015]

Shapes



[Monti et al., 2016]



Link prediction with Graph Auto-Encoders

Kipf & Welling, NIPS Bayesian Deep Learning Workshop, 2016





Further reading

Blog post Graph Convolutional Networks: http://tkipf.github.io/graph-convolutional-networks

Code on Github: http://github.com/tkipf/gcn

Kipf & Welling, Semi-Supervised Classification with Graph Convolutional Networks, ICLR 2017: <u>https://arxiv.org/abs/1609.02907</u>

Kipf & Welling, Variational Graph Auto-Encoders, NIPS BDL Workshop, 2016: https://arxiv.org/abs/1611.07308

You can get in touch with me via:

- E-Mail: T.N.Kipf@uva.nl
- Twitter: @thomaskipf
- Web: http://tkipf.github.io



