# GRAPH ATTENTION NETWORKS

Petar Velickovic´, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, Yoshua Bengio
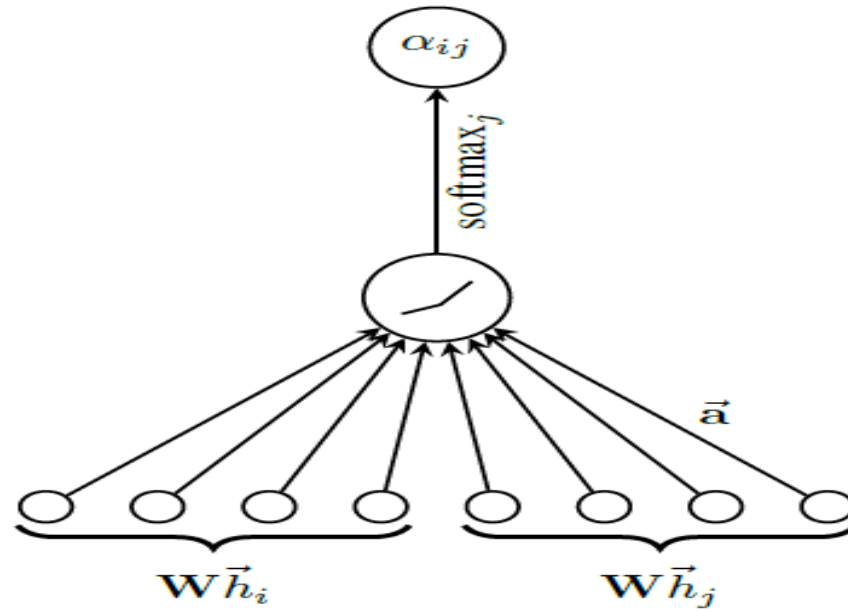
Presented by: Tawkat

# Why Graph Neural Net

- CNN can tackle problem when the underlying data representation is a grid-like structure.

- Not in the case of social networks, 3D-Mesh, biological networks

- We need to represent them with graph structure

# Why Graph Attn Net

- Traditional GNN assigns same weights to all nodes within a neighborhood, whereas GAT assigns weights to the nodes depending on their contributions.

- GAT enables better interpretability

- GAT applies attention mechanism to all the edges, so it does not depend on the global graph structure.

# Graph Attn Net



- Input node features: $h = \{h_1, h_2, ..., h_N\} \in R^F$
- Output node features: $h' = \{h'_1, h'_2, ..., h'_N\} \in R^{F'}$
- Weight Matrix: $W \in R^{F \times F'}$
- eij = $a(Wh_i, Wh_j)$; where a: $R^{2 \times F'}$ is attention coeff. that indicate the importance of node j's features to node i.
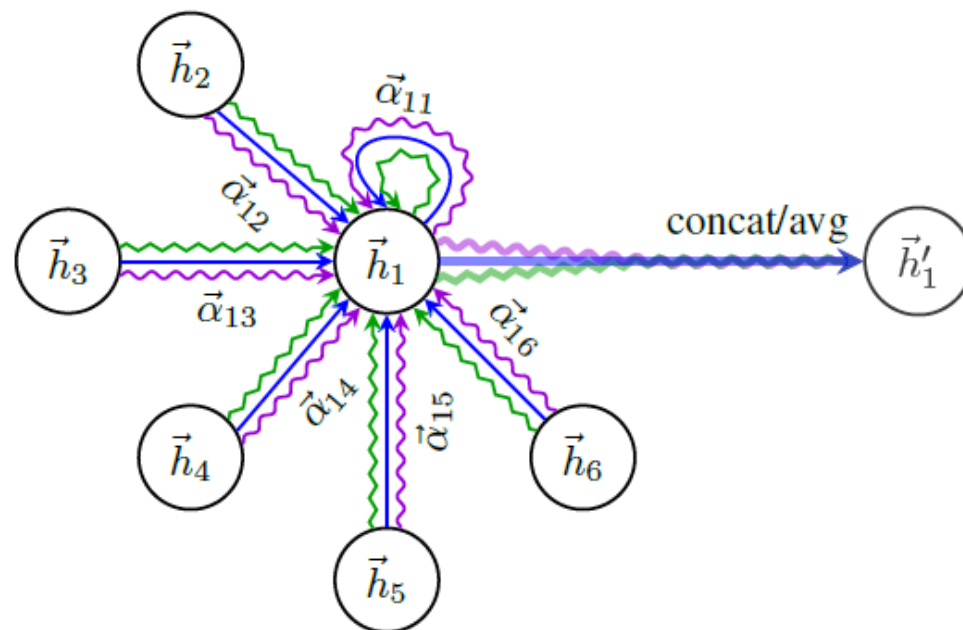
# Graph Attn Net

$$e_{ij} = a(\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_j)$$

$$\alpha_{ij} = \mathrm{softmax}_j(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}_i} \exp(e_{ik})}.$$

$$\alpha_{ij} = \frac{\exp\left(\mathrm{LeakyReLU}\left(\vec{\mathbf{a}}^T[\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_j]\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\mathrm{LeakyReLU}\left(\vec{\mathbf{a}}^T[\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_k]\right)\right)}$$

$$\vec{h}_i' = \sigma\left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W}\vec{h}_j\right).$$

# Graph Attn Net-Multi-Head



$$\vec{h}_i' = \overset{K}{\underset{k=1}{\Big\|}}\ \sigma\left(\sum_{j\in\mathcal{N}_i}\alpha_{ij}^k\mathbf{W}^k\vec{h}_j\right)$$

$$\vec{h}_i' = \sigma\left(\frac{1}{K}\sum_{k=1}^{K}\sum_{j\in\mathcal{N}_i}\alpha_{ij}^k\mathbf{W}^k\vec{h}_j\right)$$

# Experimental Results

**Transductive**

| Method | Cora | Citeseer | Pubmed |
|---|---|---|---|
| MLP | 55.1% | 46.5% | 71.4% |
| ManiReg (Belkin et al., 2006) | 59.5% | 60.1% | 70.7% |
| SemiEmb (Weston et al., 2012) | 59.0% | 59.6% | 71.7% |
| LP (Zhu et al., 2003) | 68.0% | 45.3% | 63.0% |
| DeepWalk (Perozzi et al., 2014) | 67.2% | 43.2% | 65.3% |
| ICA (Lu & Getoor, 2003) | 75.1% | 69.1% | 73.9% |
| Planetoid (Yang et al., 2016) | 75.7% | 64.7% | 77.2% |
| Chebyshev (Defferrard et al., 2016) | 81.2% | 69.8% | 74.4% |
| GCN (Kipf & Welling, 2017) | 81.5% | 70.3% | **79.0%** |
| MoNet (Monti et al., 2016) | $81.7 \pm 0.5\%$ | — | $78.8 \pm 0.3\%$ |
| GCN-64* | $81.4 \pm 0.5\%$ | $70.9 \pm 0.5\%$ | $\mathbf{79.0} \pm 0.3\%$ |
| **GAT** (ours) | $\mathbf{83.0} \pm 0.7\%$ | $\mathbf{72.5} \pm 0.7\%$ | $\mathbf{79.0} \pm 0.3\%$ |

# Experimental Results

| Method | PPI |
|---|---|
| Random | 0.396 |
| MLP | 0.422 |
| GraphSAGE-GCN (Hamilton et al., 2017) | 0.500 |
| GraphSAGE-mean (Hamilton et al., 2017) | 0.598 |
| GraphSAGE-LSTM (Hamilton et al., 2017) | 0.612 |
| GraphSAGE-pool (Hamilton et al., 2017) | 0.600 |
| GraphSAGE* | 0.768 |
| Const-GAT (ours) | $0.934 \pm 0.006$ |
| **GAT** (ours) | $\mathbf{0.973} \pm 0.002$ |

# Future Works

- Analysis on better interpretability.
- Graph classification instead of node classification.
- Including edge features instead of treating edges just as boolean variables