

FastGCN: Fast Learning with Graph Convolution Networks Via Importance Sampling

Jie Chen, Tengfei Ma, Cao Xiao
IBM Research

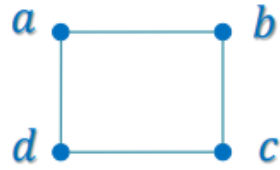
ICLR 2018

Sanghyeon Lee

26 Aug. 2020

Backgrounds

1) Property of Adjacency Matrices (A)



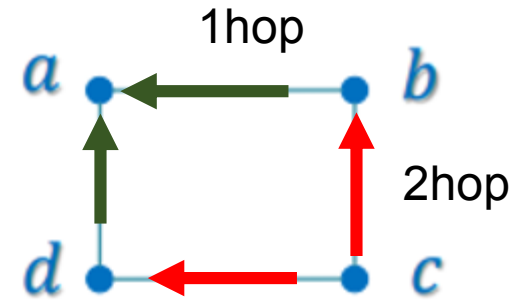
$$A = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix}$$

$$A^2 = AA = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 2 & 0 & 2 & 0 \\ 0 & 2 & 0 & 2 \\ 2 & 0 & 2 & 0 \\ 0 & 2 & 0 & 2 \end{pmatrix}$$

$$A^3 = A^2A = \begin{pmatrix} 2 & 0 & 2 & 0 \\ 0 & 2 & 0 & 2 \\ 2 & 0 & 2 & 0 \\ 0 & 2 & 0 & 2 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 4 & 0 & 4 \\ 4 & 0 & 4 & 0 \\ 0 & 4 & 0 & 4 \\ 4 & 0 & 4 & 0 \end{pmatrix}$$

2) Adjacency Matrices and Multi hop

- Aggregating feature from connected 1 hop node



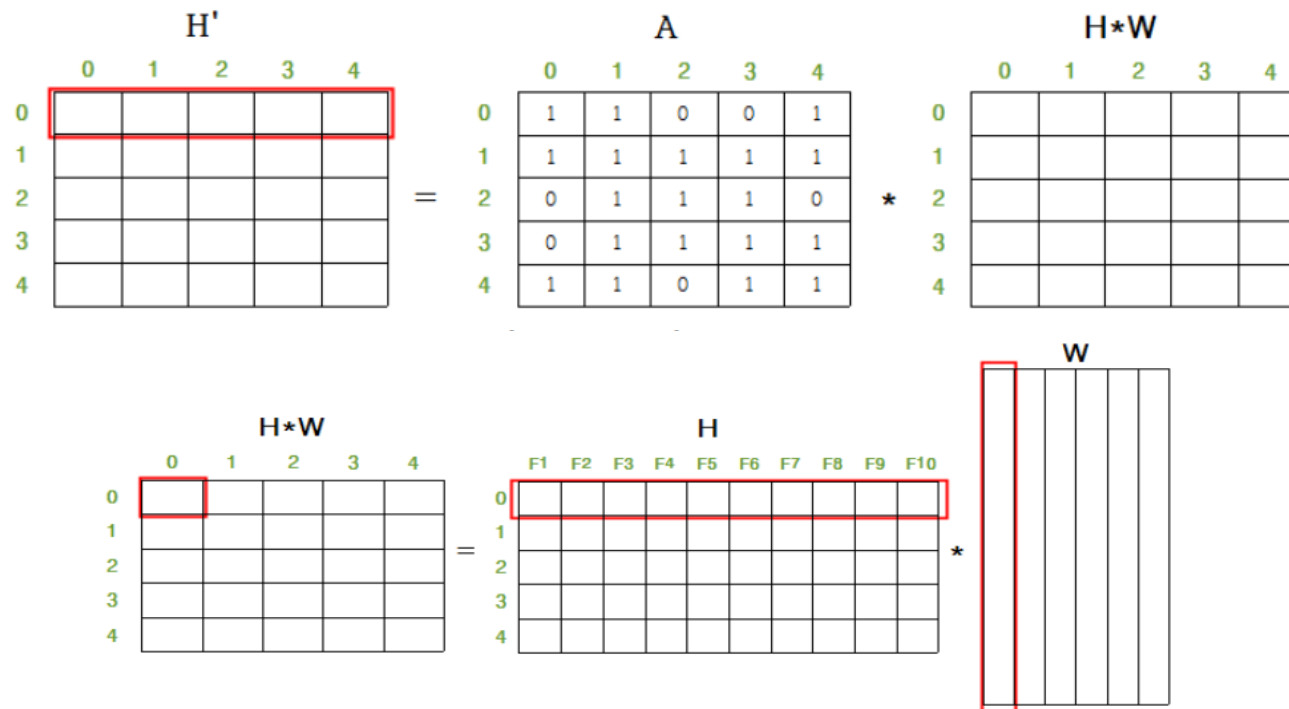
Aggeratied node feature = AH

Backgrounds

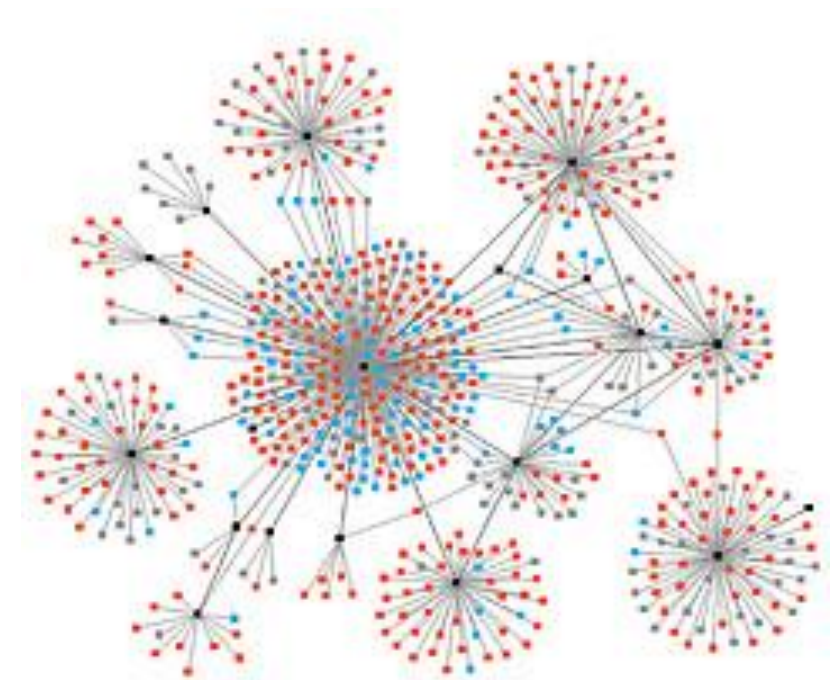
3) GCN

$$H_1^{(l+1)} = \sigma \left(H_1^{(l)} W^{(l)} + H_2^{(l)} W^{(l)} + H_3^{(l)} W^{(l)} + H_4^{(l)} W^{(l)} + b^{(l)} \right)$$

$$\rightarrow H_i^{(l+1)} = \sigma \left(\sum_{j \in N(i)} H_j^{(l)} W^{(l)} + b^{(l)} \right) \quad H^{(l+1)} = \sigma \left(A H^{(l)} W^{(l)} + b^{(l)} \right)$$



small-world network

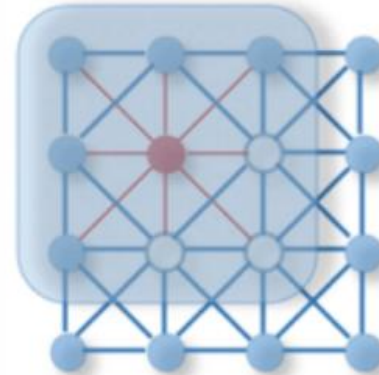
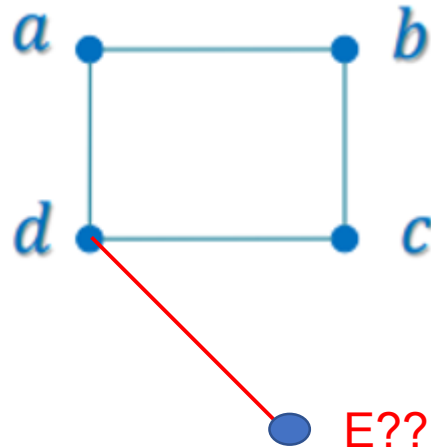


Backgrounds

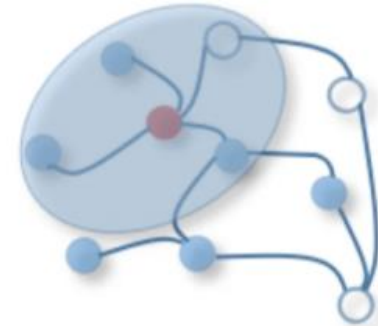
Problem of GCN:

1. For many applications, test graph may be constantly expanding with new nodes
2. Very Large Computational Cost

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



(a) 2D Convolution. Analogous to a graph, each pixel in an image is taken as a node where neighbors are determined by the filter size. The 2D convolution takes a weighted average of pixel values of the red node along with its neighbors. The neighbors of a node are ordered and have a fixed size.



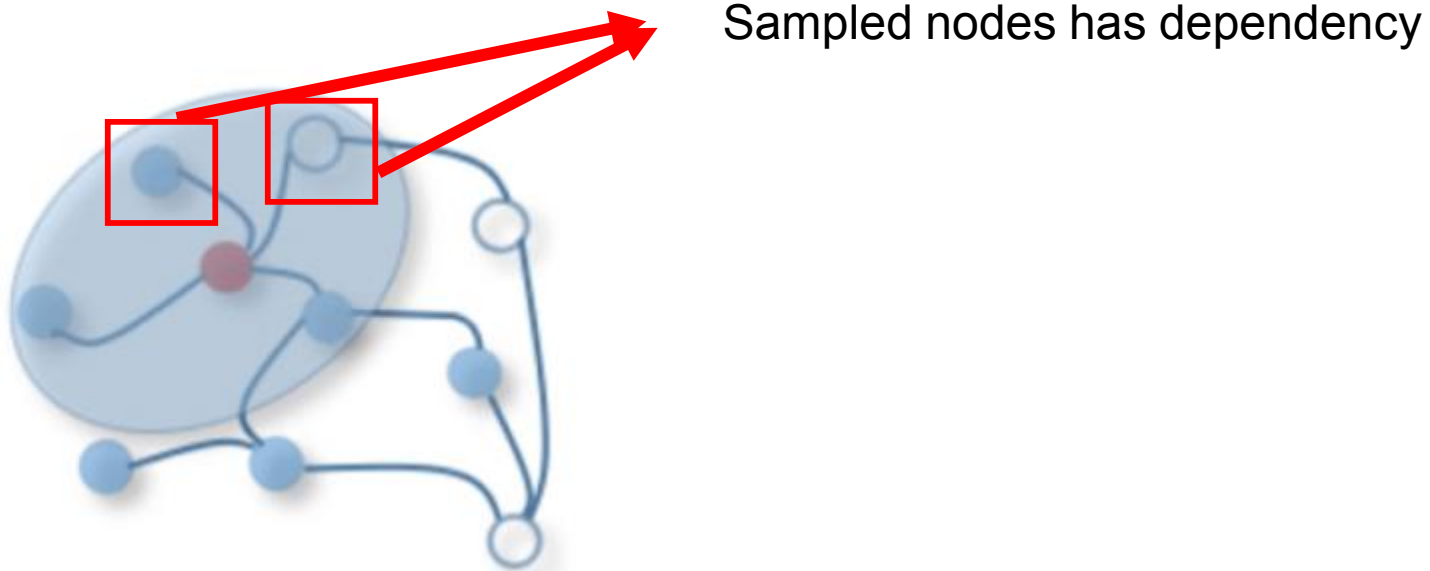
(b) Graph Convolution. To get a hidden representation of the red node, one simple solution of graph convolution operation takes the average value of node features of the red node along with its neighbors. Different from image data, the neighbors of a node are unordered and variable in size.

Fig. 1: 2D Convolution vs. Graph Convolution.

Backgrounds

Problem of Batch sampling in GCN:

3. Nodes has high dependency \rightarrow SGD assumes that all samples are iid



Method

Training and Inference through Sampling

$$\tilde{H}^{(l+1)} = \hat{A}H^{(l)}W^{(l)}, \quad H^{(l+1)} = \sigma(\tilde{H}^{(l+1)}), \quad l = 0, \dots, M-1,$$

Let function $h, s.t$

$$\tilde{h}^{(l+1)}(v) = \int \hat{A}(v, u)h^{(l)}(u)W^{(l)} dP(u), \quad h^{(l+1)}(v) = \sigma(\tilde{h}^{(l+1)}(v)), \quad l = 0, \dots, M-1,$$

By Monte Carlo Sampling

$$\tilde{h}_{t_{l+1}}^{(l+1)}(v) := \frac{1}{t_l} \sum_{j=1}^{t_l} \hat{A}(v, u_j^{(l)})h_{t_l}^{(l)}(u_j^{(l)})W^{(l)}, \quad h_{t_{l+1}}^{(l+1)}(v) := \sigma(\tilde{h}_{t_{l+1}}^{(l+1)}(v)), \quad u \sim P$$

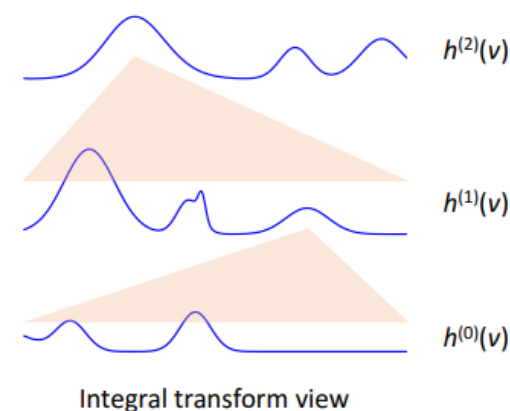
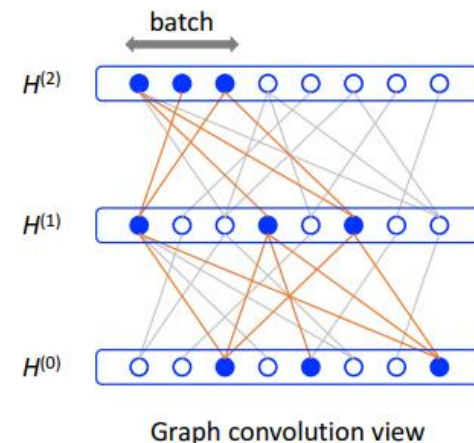
Vertex Sampling

$$H^{(l+1)}(v, :) = \sigma \left(\frac{n}{t_l} \sum_{j=1}^{t_l} \hat{A}(v, u_j^{(l)})H^{(l)}(u_j^{(l)}, :)W^{(l)} \right), \quad l = 0, \dots, M-1.$$

Batch Sampling

$$L_{t_0, t_1, \dots, t_M} := \frac{1}{t_M} \sum_{i=1}^{t_M} g(h_{t_M}^{(M)}(u_i^{(M)}))$$

They assumes that all embedded nodes follows iid → This assumption is too strong ...



Method

*** Adaptive Sampling Toward Fast Graph Representation Learning (NIPS 2018)**

Vertex Sampling

$$H^{(l+1)}(v, :) = \sigma \left(\frac{n}{t_l} \sum_{j=1}^{t_l} \hat{A}(v, u_j^{(l)}) H^{(l)}(u_j^{(l)}, :) W^{(l)} \right), \quad l = 0, \dots, M-1.$$

$$pf) \sigma(\Sigma_j A(v_i, u_j) h(u_j) W) = \sigma \left(\frac{n}{1} * \frac{1}{n} \dots \right) = \sigma \left(n * E_{p(u_j|v_i)}(h(u_j) W) \right)$$

$$s.t \ p(u_j|v_i) = \frac{A(v_i, u_j)}{n}$$

$$\Rightarrow H = \sigma \left(n * \frac{1}{t_l} \Sigma \dots \right)$$

Method

Training and Inference through Sampling

Algorithm 1 FastGCN batched training (one epoch)

- 1: **for** each batch **do**
- 2: For each layer l , sample uniformly t_l vertices $u_1^{(l)}, \dots, u_{t_l}^{(l)}$
- 3: **for** each layer l **do** ▷ Compute batch gradient ∇L_{batch}
- 4: If v is sampled in the next layer,

$$\nabla \tilde{H}^{(l+1)}(v, :) \leftarrow \frac{n}{t_l} \sum_{j=1}^{t_l} \hat{A}(v, u_j^{(l)}) \nabla \left\{ H^{(l)}(u_j^{(l)}, :) W^{(l)} \right\}$$

- 5: **end for**
 - 6: $W \leftarrow W - \eta \nabla L_{\text{batch}}$ ▷ SGD step
 - 7: **end for**
-

Method

Variance Reduction by importance sampling

	Function	Samples	Num. samples
Layer $l + 1$; random variable v	$\tilde{h}_{t_{l+1}}^{(l+1)}(v) \rightarrow y(v)$	$u_i^{(l+1)} \rightarrow v_i$	$t_{l+1} \rightarrow s$
Layer l ; random variable u	$h_{t_l}^{(l)}(u)W^{(l)} \rightarrow x(u)$	$u_j^{(l)} \rightarrow u_j$	$t_l \rightarrow t$

Under the joint distribution of v and u , the aforementioned sample average is

$$G := \frac{1}{s} \sum_{i=1}^s y(v_i) = \frac{1}{s} \sum_{i=1}^s \left(\frac{1}{t} \sum_{j=1}^t \hat{A}(v_i, u_j) x(u_j) \right).$$

$$\text{Var}\{G\} = R + \frac{1}{st} \iint \hat{A}(v, u)^2 x(u)^2 dP(u) dP(v), \quad (6)$$

where

$$R = \frac{1}{s} \left(1 - \frac{1}{t} \right) \int e(v)^2 dP(v) - \frac{1}{s} \left(\int e(v) dP(v) \right)^2 \quad \text{and} \quad e(v) = \int \hat{A}(v, u) x(u) dP(u).$$

Method

Variance Reduction by importance sampling

$$\text{Importance Sampling : } E_p(f(x)) = E_q\left(\frac{p}{q}f(x)\right)$$

Algorithm 2 FastGCN batched training (one epoch), improved version

- 1: For each vertex u , compute sampling probability $q(u) \propto \|\hat{A}(:, u)\|^2$
- 2: **for** each batch **do**
- 3: For each layer l , sample t_l vertices $u_1^{(l)}, \dots, u_{t_l}^{(l)}$ according to distribution q
- 4: **for** each layer l **do** \triangleright Compute batch gradient ∇L_{batch}
- 5: If v is sampled in the next layer,

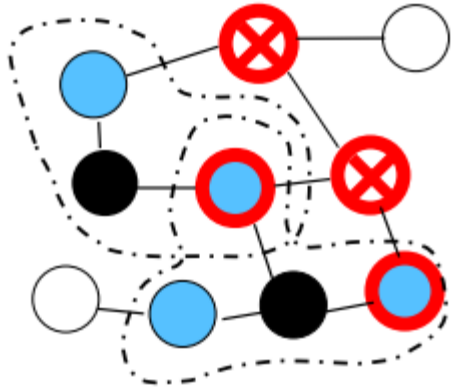
$$\nabla \tilde{H}^{(l+1)}(v, :) \leftarrow \frac{1}{t_l} \sum_{j=1}^{t_l} \frac{\hat{A}(v, u_j^{(l)})}{q(u_j^{(l)})} \nabla \left\{ H^{(l)}(u_j^{(l)}, :) W^{(l)} \right\}$$

- 6: **end for**
 - 7: $W \leftarrow W - \eta \nabla L_{\text{batch}}$ \triangleright SGD step
 - 8: **end for**
-

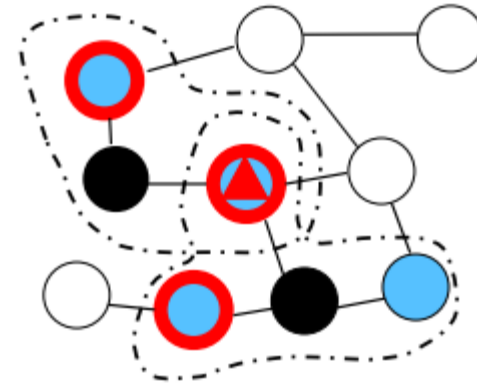
Discussion

Problem of FastGCN

1.



**Layer-wise Importance Sampling
(FastGCN)**



**Node-wise Neighbor Sampling
(GraphSAGE)**

2. Batch sampling has less theoretical background

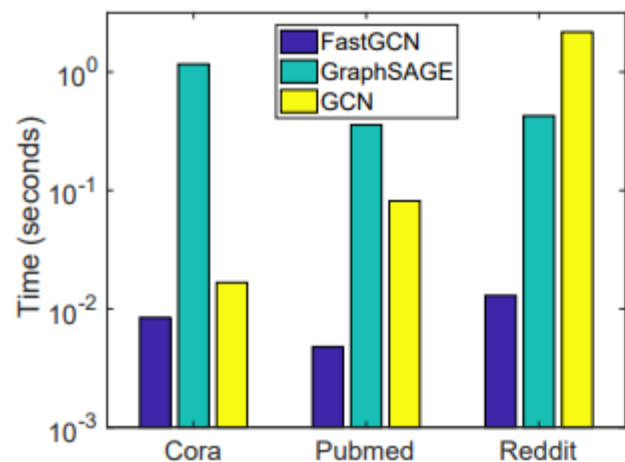
Experiments

- 1 Classifying research topics using Cora citation data set
- 2 Categorizing academic papers with the Pubmed database
- 3 Predicting the community structure of a social network modeled with Reddit posts

Table 1: Dataset Statistics

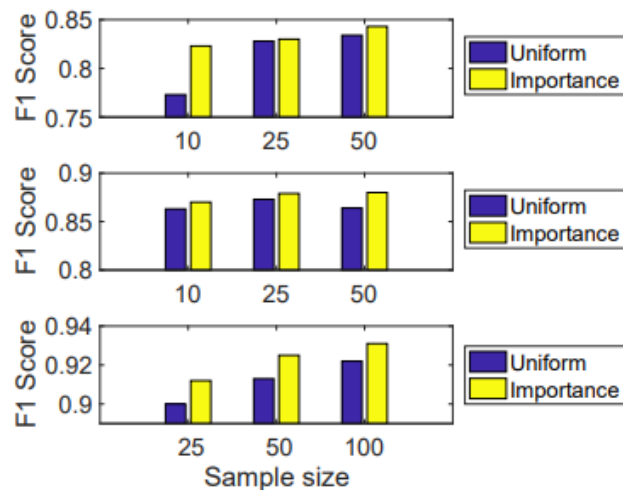
Dataset	Nodes	Edges	Classes	Features	Training/Validation/Test
Cora	2,708	5,429	7	1,433	1,208/500/1,000
Pubmed	19,717	44,338	3	500	18,217/500/1,000
Reddit	232,965	11,606,919	41	602	152,410/23,699/55,334

Experiments



Micro F1 Score			
	Cora	Pubmed	Reddit
FastGCN	0.850	0.880	0.937
GraphSAGE-GCN	0.829	0.849	0.923
GraphSAGE-mean	0.822	0.888	0.946
GCN (batched)	0.851	0.867	0.930
GCN (original)	0.865	0.875	NA

t_1	Sampling		Precompute	
	Time	F1	Time	F1
5	0.737	0.859	0.139	0.849
10	0.755	0.863	0.141	0.870
25	0.760	0.873	0.144	0.879
50	0.774	0.864	0.142	0.880



Thank you