

# NetGAN without GAN: From Random Walks to Low-Rank Approximations

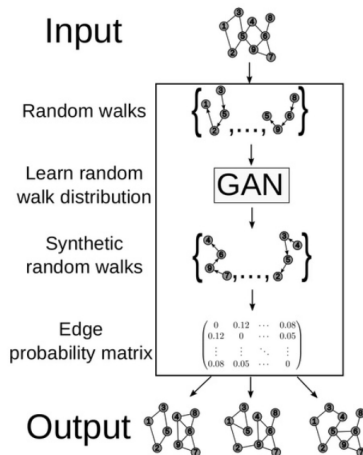
Numerical Linear Algebra course final project

Evgeny Lagutin, Pavel Burnyshev, Valentin Samokhin, Daniil Moskovskiy

December 19, 2020

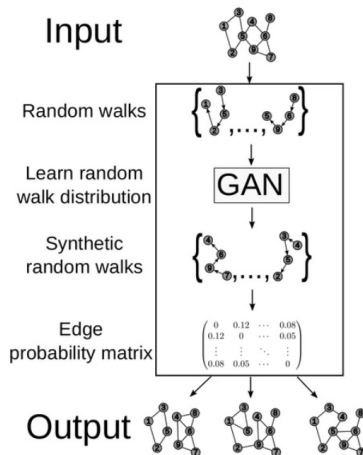
NetGAN aims to reconstruct graphs similar to the input one. NetGAN algorithm in short:

- 1 Sample a set of random walks from the input graph to train the GAN
- 2 Learn random walk distribution
- 3 Sample a set of random walks from the GAN
- 4 Construct edge probability matrix from the produced set



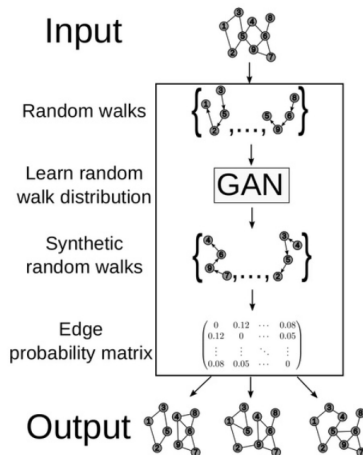
NetGAN aims to reconstruct graphs similar to the input one. NetGAN algorithm in short:

- 1 Sample a set of random walks from the input graph to train the GAN
- 2 Learn random walk distribution
- 3 Sample a set of random walks from the GAN
- 4 Construct edge probability matrix from the produced set



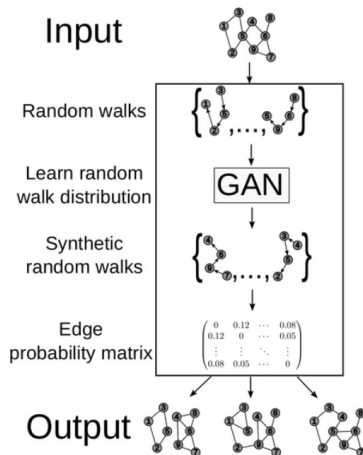
NetGAN aims to reconstruct graphs similar to the input one. NetGAN algorithm in short:

- 1 Sample a set of random walks from the input graph to train the GAN
- 2 Learn random walk distribution
- 3 Sample a set of random walks from the GAN
- 4 Construct edge probability matrix from the produced set



NetGAN aims to reconstruct graphs similar to the input one. NetGAN algorithm in short:

- 1 Sample a set of random walks from the input graph to train the GAN
- 2 Learn random walk distribution
- 3 Sample a set of random walks from the GAN
- 4 Construct edge probability matrix from the produced set

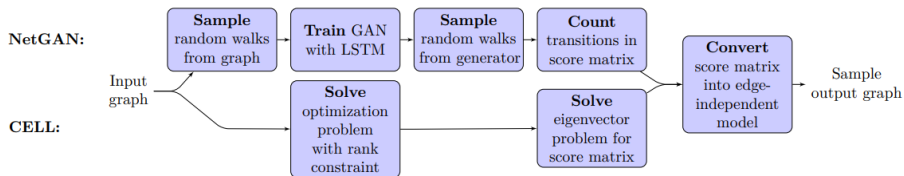


**Main goal:** Learn random walk distribution from a training set based on unbiased random walks  $\{v_0^i, \dots, v_T^i\}_i$  over input graph. It's supposed to be connected and non-bipartite.

**Generator** generates sequences  $\{\omega_0^i, \dots, \omega_T^i\}_i$  of "synthetic" random walks while **discriminator** tries to distinguish "synthetic" random walks from the real ones.

**Core model** is LSTM with Wasserstein loss (therefore, no long-term dependencies due to LSTM usage, and comparable slow training). Update equations for **synthetic** random walk are given by:

$$(m_{t+1}, p_{t+1}) = f_{\theta}(m_t, v_t)$$
$$v_{t+1} \sim \text{Cat}(\sigma(p_{t+1}))$$



- Replaces the original problem to a low-rank approximation with respect to the Kullback-Leibler divergence between transition matrices (requires neither a GAN nor any sampling).
- Increases the speed of training and is more transparent due to simpler structure
- Evades random walk sampling part which was (will be shown further) very computationally expensive and inefficient.

The learnable projection matrices are given by  $W_{\text{down}} \in \mathbb{R}^{N \times H}$  and  $W_{\text{up}} \in \mathbb{R}^{H \times N}$  with  $H \ll N$ . Given the current node  $v_t$  as a one-hot vector, the generator  $f_\theta$  can be written as

$$(p_{t+1}, m_{t+1}) = f_\theta(m_t, v_t) = g_\theta\left(m_t, v_t^\top W_{\text{down}}\right) W_{\text{up}}$$

where  $g_\theta : \mathbb{R}^H \rightarrow \mathbb{R}^H$  is the part of  $f_\theta$  that operates on the low-dimensional space. We collect the row vectors  $g_\theta(m_t, v_t^\top W_{\text{down}})$  in a matrix  $\widetilde{W}_{\text{down}}(m_t) \in \mathbb{R}^{N \times H}$  and define the product  $W(m_t) := \widetilde{W}_{\text{down}}(m_t) W_{\text{up}} \in \mathbb{R}^{N \times N}$  to obtain

$$p_{t+1} = v_t^\top \widetilde{W}_{\text{down}}(m_t) W_{\text{up}} = v_t^\top W(m_t)$$

Therefore,  $W(m_t)$  simply serves as logit transition matrix for the r.w. Because of the factorization that defines  $W(m_t)$ , its rank is at most  $H$



The authors propose dropping the LSTM with its dependency on the memory state  $m_t$ ; this is justified by the Markov property of unbiased r.w.:

$$p_{t+1} = v_t^\top W = g_\theta \left( v_t^\top W_{\text{down}} \right) W_{\text{up}}$$

Authors suggest to drop the structural restriction imposed by  $g_\theta$ , setting  $W = W_{\text{down}} W_{\text{up}}$  and update rule

$$p_{t+1} = v_t^\top W_{\text{down}} W_{\text{up}}$$

But we may leave the  $g_\theta$  part by rewriting as ( $v_\theta$  is one-hot vector):

$$p_{t+1} = v_t^\top g_\theta(W_{\text{down}}) W_{\text{up}}$$

Therefore  $v_{t+1} \sim \text{Cat}(\sigma(W_{v_t}))$ .

Now training the GAN is equivalent to learning the r.w. transition matrix from the family  $\mathcal{P} = \{ \sigma_{\text{rows}}(W) \in \mathbb{R}^{N \times N} : W \in \mathbb{R}^{N \times N}, \text{rank}(W) \leq H \}$

Consider a graph with number of nodes  $N$ , set of random walks  $\mathcal{R}$ , low rank  $H$  row-wise softmax  $\sigma_{\text{rows}}$

Learn random walk distribution by **learning random walk transition matrix** from parametric family:

$$\mathcal{P} = \{\sigma_{\text{rows}}(W) \in \mathbb{R}^{N \times N} : W \in \mathbb{R}^{N \times N}, \text{rank}(W) \leq H\}$$

Solve an optimization problem: 
$$\begin{cases} - \sum_{(i,j) \in \mathcal{R}} \log \sigma_{\text{rows}}(W)_{i,j} \rightarrow \min_{W \in \mathbb{R}^{N \times N}} \\ \text{s.t.: rank}(W) \leq H \end{cases}$$

Latent variable  $z$  plays the subordinate role of choosing the first node.

NetGAN samples a huge amount of random walks from generator that learned the distribution from the input graph

E.g., for CORA-ML graph (2810 nodes, 7891 edges) NetGAN samples that much random walks, that every edge of the graph is processed for 14000 times (!).

Given length of random walks  $T$ , number of random walks  $n$ , score matrix  $S \in \mathbb{R}^{N \times N}$ , the normalized score matrix converges:

$$\frac{S}{nT} \xrightarrow[n, T \rightarrow +\infty]{a.s.} \text{diag}(\pi)P,$$

where  $P$  is corresponding random walk transition matrix and  $\pi$  is stationary distribution.

**Learning step** with adjacency matrix  $A$

$$\begin{cases} - \sum_{(i,j) \in \mathcal{R}} \log \sigma_{\text{rows}}(W)_{i,j} \rightarrow \min_{W \in \mathbb{R}^{N \times N}} \\ \text{s.t.: rank}(W) \leq H \end{cases} \quad \begin{cases} - \sum_{k,l} A_{k,l} \log \sigma_{\text{rows}}(W)_{k,l} \rightarrow \min_{W \in \mathbb{R}^{N \times N}} \\ \text{s.t.: rank}(W) \leq H \end{cases}$$

**Reconstruction step** (with synthetic transition matrix  $P^*$ ):

- 1 Find  $\pi^*$  s.t.  $\pi^{*\top} = \pi^{*\top} P^*$
- 2 Take  $S := \text{diag}(\pi^*) P^*$

**No sampling during both steps!!!**

---

## Algorithm 1: CELL algorithm

---

**Input:** adjacency matrix  $A \in \{0, 1\}^{N \times N}$ , rank  $H \ll N$

**Output:** matrix of edge probabilities  $A^\dagger \in [0, 1]^{N \times N}$

- 1 Solve optimization problem for  $W^*$  (shown on the previous slide)
  - 2 Compute transition matrix:  $P^* \leftarrow \sigma_{\text{rows}}(W^*)$
  - 3 Solve eigenvalue problem:  $\pi^{*T} P^* = \pi^{*T}$  for  $\pi^*$
  - 4 Compute score matrix:  $S \leftarrow \text{diag}(\pi^*) P^*$
  - 5 Convert score matrix  $S$  to edge-independent model  $A^\dagger$ :  
 $S^\dagger \leftarrow \max \{S, S^T\}$ ;  $A^\dagger \leftarrow S^\dagger / \sum_{i,j} S_{i,j}^\dagger$ , so  
$$A_{i,j}^\dagger = \max \{S_{i,j}\} / \sum_{k',l'=1}^N \max \{S_{k,l}, S_{k',l'}\}$$
  - 6 Return  $A^\dagger$
-

Assortativity	$\frac{\text{cov}X,Y}{\sigma_X\sigma_Y}$
Power Law Exponent	$1 + n \left( \sum_{v \in V} \log \frac{d(v)}{d_{\min}} \right)^{-1}$
Relative edge distr. entropy	$-\frac{1}{\log N} \sum_{v \in V} \frac{d(v)}{2m} \log \frac{d(v)}{2m}$
Gini coeff.	$\frac{2 \sum_{i=1}^N i \hat{d}_i}{\sum_{i=1}^N \hat{d}_i} - \frac{N+1}{N}$
Character. path length	$\frac{1}{N(N-1)} \sum_{u,v} d(u,v)$
Spectral gap	$\lambda_1(L)$
Motif count	--

Table: Graph statistics used

The performance is measured by the ROC-AUC score, applied to the score matrix evaluated at the edges in question (test edges)

- replication of the proposed approach and results
- extending parametric family of low rank matrices  $W$  with decomposition on non-linear transformations (implementing  $W$  as output of fully-connected auto-encoder with  $I_N$  as input)
- experimenting with low-rank approximations of  $W$  with SVD by explicitly constraining  $U, V$  to be close to orthogonal during gradient descent

Method	$d_{\max}$	Assort.	Triangle count	Square count	Power Law Exp	Clustering Coeff.	ROC-AUC
Original Graph	246	-0.077	5247	34507	1.77	0.114	1
CELL	213	-0.066	1235	5936	<b>1.81</b>	0.045	0.94
SVD	196	<b>-0.079</b>	1264	6813	1.83	0.046	0.95
nonlinear-CELL	<b>245</b>	-0.074	<b>1791</b>	<b>10122</b>	1.86	<b>0.054</b>	<b>0.96</b>

Figure: Results on CORA-ML graph (2810 nodes, 7891 edges) using edge overlap criterion.

Method	$d_{\max}$	Assort.	Triangle count	Square count	Power Law Exp	Clustering Coeff.	ROC-AUC
Original Graph	99	0.060	1083	5977	2.068	0.125	1
CELL	55	-0.060	138	413	<b>2.183</b>	0.031	0.88
SVD	60	-0.084	165	501	2.228	0.033	<b>0.92</b>
nonlinear-CELL	<b>96</b>	<b>0.017</b>	<b>441</b>	<b>3041</b>	2.330	<b>0.062</b>	0.88

Figure: Results on Citeseer graph (2110 nodes, 3668 edges) using edge overlap criterion.



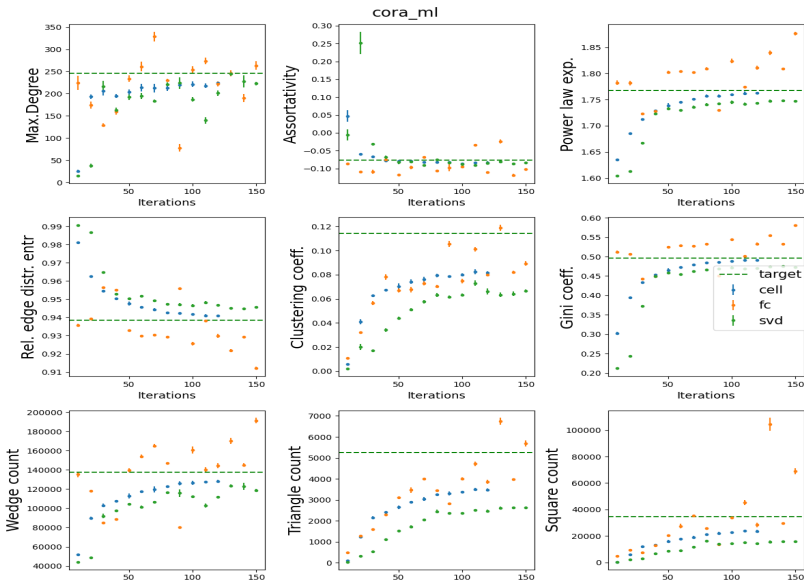
Method	$d_{\max}$	Assort.	Triangle count	Square count	Power Law Exp	Clustering Coeff.	ROC-AUC
Original Graph	351	-0.184	101043	5171257	1.414	0.226	1
CELL	260	<b>-0.225</b>	40880	1576069	1.407	0.144	<b>0.96</b>
SVD	234	-0.231	<b>43161</b>	<b>1725619</b>	<b>1.412</b>	<b>0.149</b>	<b>0.96</b>
nonlinear-CELL	<b>276</b>	-0.229	42507	1696934	1.419	0.148	0.95

Figure: Results on PolBlogs graph (1222 nodes, 16779 edges) using edge overlap criterion.

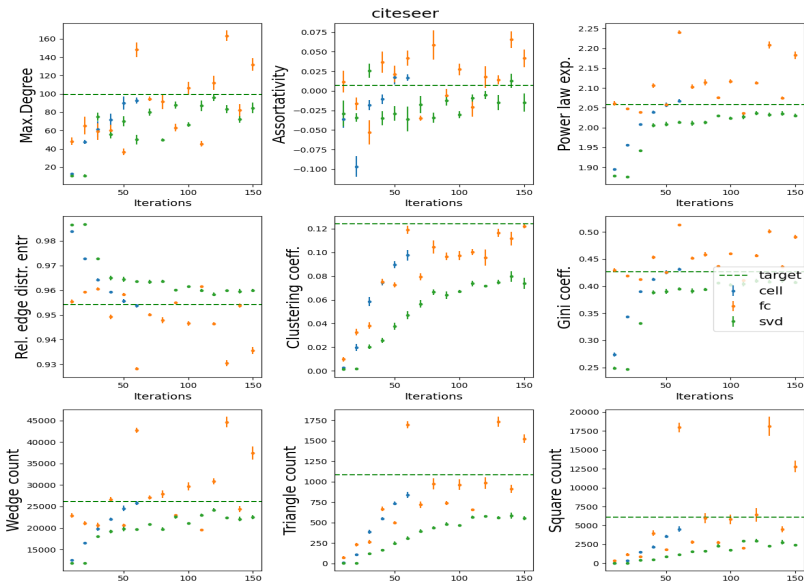
Method	$d_{\max}$	Assort.	Triangle count	Square count	Power Law Exp	Clustering Coeff.	ROC-AUC
Original Graph	308	-0.090	33	369	3.493	$0.0008 \approx 0$	1
CELL	<b>248</b>	-0.142	0	11	<b>3.942</b>	0	0.69
SVD	234	<b>-0.137</b>	0	10	4.149	0	0.77
nonlinear-CELL	242	-0.142	0	<b>13</b>	4.269	0	<b>0.84</b>

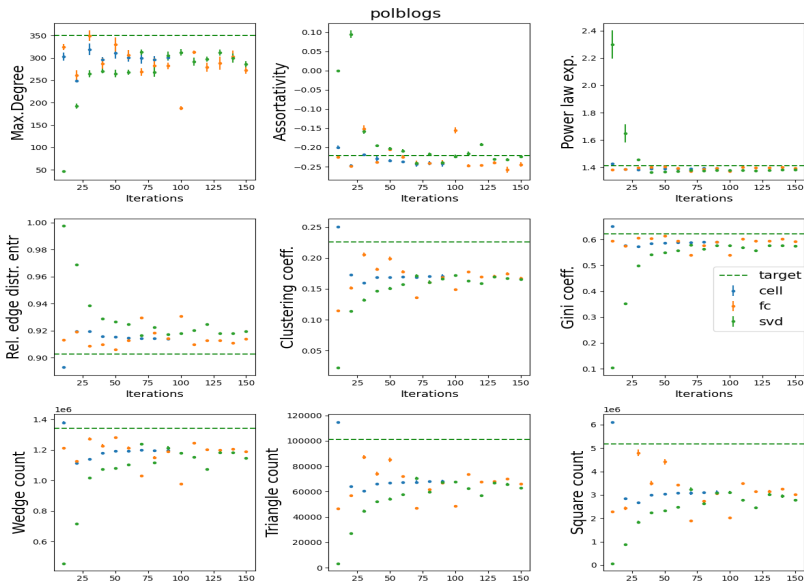
Figure: Results on RT-GOP graph (4687 nodes, 5529 edges) using edge overlap criterion.

# Experiments: plots - CORAML

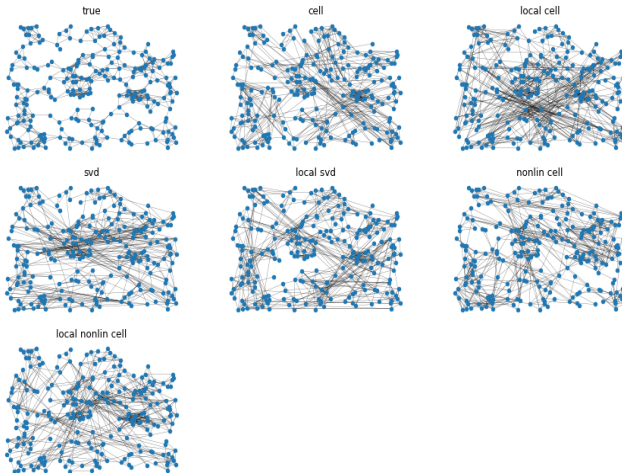


# Experiments: plots - Citeseer





# Experiments: $\varepsilon$ -neighbourhood graph



■ Thank you for attention!