

# A Discussion on “Graphite: Iterative Generative Modeling of Graphs”

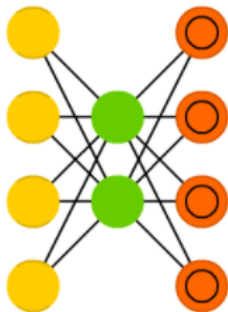
Discussant: Roman Kouznetsov

October 19, 2020

# Review of VAE

VAE was introduced as a way to turn the bottleneck of autoencoders into a generative process from a latent space.

Auto Encoder (AE)



Variational AE (VAE)

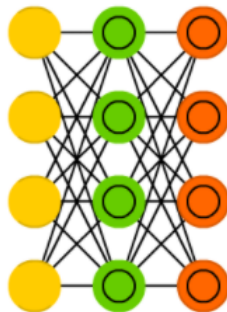


Figure: AE v. VAE

# Graphs

- ▶ Graphite works with undirected, unweighted graphs so that is the kind of graph I will introduce here.
- ▶ Graphs are objects that contain nodes and edges that connect them.
- ▶ Unweighted graphs means that all edges are treated identically storing the same information.
- ▶ Undirected graphs means that edges are bidirectional (i.e. messages can be passed both ways). This means that the adjacency matrix is symmetrical (more on this later).
- ▶ Graphs can have node labels, node features, and edge features that better describe the structure of a graph.

# Graph Examples

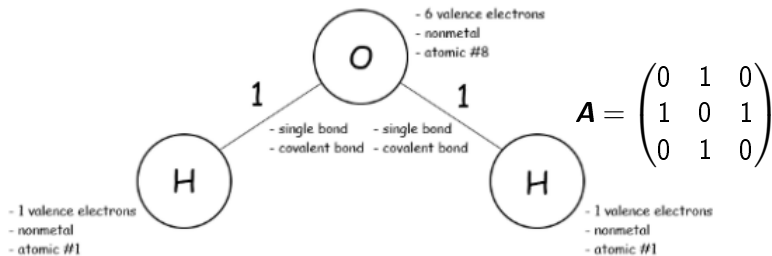


Figure: A Water Molecule Graph w/ Corresponding Adjacency Matrix

Image Credit: <https://towardsdatascience.com/graph-theory-and-deep-learning-know-hows-6556b0e9891b>

# Some Graph Notation

- ▶  $G = (V, E)$ : undirected graph with nodes  $V$  and edges  $E$
- ▶  $X \in \mathbb{R}^{n \times m}$ : optional feature matrix where each node has an  $m$ -dimensional attribute
- ▶  $A \in \mathbb{R}^{n \times n}$ : adjacency matrix s.t.  $A_{i,j}$  represents the weight of the edge between node  $i$  and  $j$
- ▶  $H = GNN_{\langle W, B \rangle}(A, X)$ : graph neural network of adjacency matrix  $A$  and feature matrix  $X$  on learnable weight and bias parameters  $W$  and  $B$  (where  $W$  and  $B$  contain all of the weights and biases across layers)

**Takeaway:** A graph is defined by its adjacency and feature matrices. This will be an important concept later.

# Introduction to Graph Neural Networks (GNNs)

$$\mathbf{H}^{(l)} \leftarrow \eta_l \left( \mathbf{B}_l + \sum_{f \in \mathcal{F}_l} f(\mathbf{A}) \mathbf{H}^{(l-1)} \mathbf{W}_l \right)$$

$\mathbf{B}_l$ : bias at current layer

$\mathbf{A}$ : adjacency matrix

$f \in \mathcal{F}_l$ : a set of linear transformations that acts on  $\mathbf{A}$

$\mathbf{H}^{(l-1)}$ : activations of the the previous layer

$\mathbf{W}_l$ : weights at current layer

# Introduction to Graphite

- ▶ **Research Objective:** GNNs do the simple encoding while Graphite focuses on the challenges of creating the decoder.
- ▶ **Previous Work:** Previous work in this field has been done, but none have used iterative graph building with learnable parameters.
- ▶ **Laymen's Terms:** Graphite has a superior performance because of its parameterization.

# Describing the Graphite Model

Graphite introduces:

- ▶ Latent Variable Vectors:  $\mathbf{Z}_i \in \mathbb{R}^k$
- ▶ Individual Node Feature Vectors:  $\mathbf{X}_i \in \mathbb{R}^m$

These matrices will prove to be instrumental to how an iterative generative model of graphs.



## Graphite Latent Variable Model

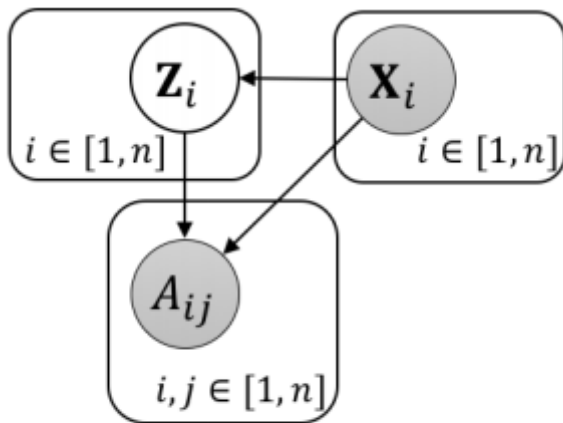


Figure 1. Latent variable model for Graphite. Observed evidence variables in gray.

Figure: Graphite Latent Variable Model

# Graphite Achievements

Graphite has been able to achieve a superior performance to its predecessor in the following problems:

- ▶ Reconstruction & Density Estimation
- ▶ Link Prediction
- ▶ Semi-Supervised Node Classification

# Encoder

Let's compare Graphite to a standard VAE.

$$(\mu, \log \sigma) = \text{EncoderNeuralNet}_{\phi}(\mathbf{x})$$

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = N(\mathbf{z}; \mu, \text{diag}\sigma)$$

$$\mu, \sigma = \text{GNN}_{\phi}(\mathbf{A}, \mathbf{X})$$

where  $\mu$  and  $\sigma$  denote the vector of means and standard deviations for the variational marginals  $\{q_{\phi_i}(\mathbf{Z}_i|\mathbf{A}, \mathbf{X})\}$

Figure: Comparing Encoders

Image Credit: Ismael Mendoza (from this reading group :-)) (left) and *Graphite: Iterative Generative Modeling of Graphs* (right)

# Decoder

*In VAE,  $\mathbf{x}$  and  $\mathbf{z}$  are vectors and so message passing between them can be established with simple NNs. However, in the context of Graphite, our input for learning  $p_{\theta}(\mathbf{x}|\mathbf{z})$  is a graph but a forward pass through the encoder yields a latent matrix. How do we create a decoder that actually has learnable parameters? **This question is the motivation for Graphite.***

$$\begin{aligned}(\mu_I, \sigma_I) &= \text{DecoderNN}_{\theta}(\mathbf{z}) & \hat{\mathbf{A}} &= \frac{\mathbf{Z}\mathbf{Z}^T}{\|\mathbf{Z}\|^2} + \mathbf{1}\mathbf{1}^T, \\ p_{\theta}(\mathbf{x}|\mathbf{z}) &= \mathcal{N}(\mathbf{x}; \mu_I, \text{diag}(\sigma_I)) & \mathbf{Z}^* &= \text{GNN}_{\theta}(\hat{\mathbf{A}}, [\mathbf{Z}|\mathbf{X}])\end{aligned}$$

Figure: Comparing Decoders

Image Credit: Ismael Mendoza (from this reading group :-)) (left) and *Graphite: Iterative Generative Modeling of Graphs* (right)

# Derivation of ELBO for GNNs

$$\begin{aligned}\log p_{\theta}(\mathbf{A}|\mathbf{X}) &= E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})}[\log p_{\theta}(\mathbf{A}|\mathbf{X})] \\&= E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{p_{\theta}(\mathbf{A}, \mathbf{Z}|\mathbf{X})}{p_{\theta}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right] \\&= E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{p_{\theta}(\mathbf{A}, \mathbf{Z}|\mathbf{X}) q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})}{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X}) p_{\theta}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right] \\&= E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{p_{\theta}(\mathbf{A}, \mathbf{Z}|\mathbf{X})}{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} + \log \frac{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})}{p_{\theta}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right] \\&= E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{p_{\theta}(\mathbf{A}, \mathbf{Z}|\mathbf{X})}{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right] + E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})}{p_{\theta}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right] \\&\geq E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{p_{\theta}(\mathbf{A}, \mathbf{Z}|\mathbf{X})}{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right]\end{aligned}$$

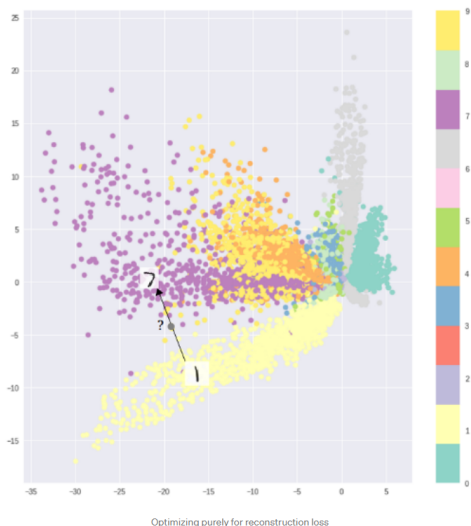
# Interpretation of ELBO

$$E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{p_{\theta}(\mathbf{A},\mathbf{Z}|\mathbf{X})}{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right] = \log p_{\theta}(\mathbf{A}|\mathbf{X}) - E_{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{q_{\phi}(\mathbf{Z}|\mathbf{A},\mathbf{X})}{p_{\theta}(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right]$$

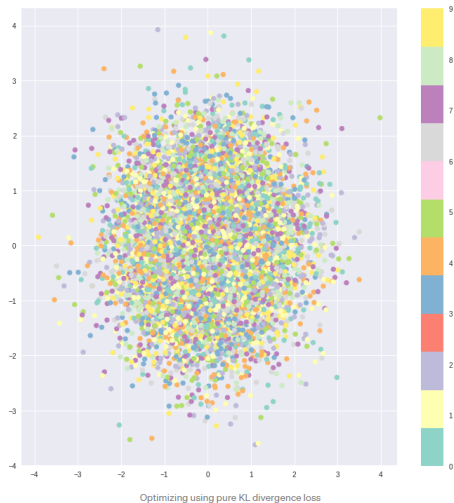
1. Learn correct adjacency matrix.
2. Make sure the learned posterior distribution is close to the (unknown) true posterior distribution.

**These 2 terms working in tandem is critical to understanding why ELBO works as such a great optimizer.**

# Showcase of Optimization on $E(\log p_{\theta}(A|Z, X))$



# Showcase of Optimization on $\text{KL}(q_\phi(z|x)||p_\theta(z|x))$



Classification X

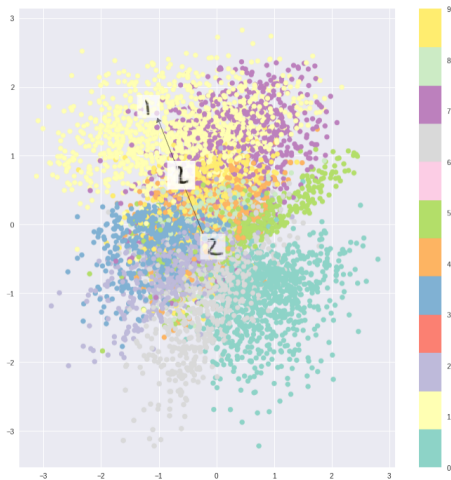
Complete Support ✓

Image Credit: <https://towardsdatascience.com/>

intuitively-understanding-variational-autoencoders-1bfe67eb5daf



# Showcase of Optimization on ELBO



Optimizing using both reconstruction loss and KL divergence loss

Classification ✓

Complete Support ✓

Image Credit: <https://towardsdatascience.com/>

intuitively-understanding-variational-autoencoders-1bfe67eb5daf

# Computation of ELBO

$$\begin{aligned} ELBO &= E_{q_\phi(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{p_\theta(\mathbf{A},\mathbf{Z}|\mathbf{X})}{q_\phi(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right] = \\ &E_{q_\phi(\mathbf{Z}|\mathbf{A},\mathbf{X})} \left[ \log \frac{p_\theta(\mathbf{A}|\mathbf{Z},\mathbf{X})p_\theta(\mathbf{Z}|\mathbf{X})}{q_\phi(\mathbf{Z}|\mathbf{A},\mathbf{X})} \right] \end{aligned}$$

Therefore, we can simply sample from the prior  $p_\theta(\mathbf{Z}|\mathbf{X})$ , our variational marginal  $q_\phi(\mathbf{Z}|\mathbf{A},\mathbf{X})$ , and our reconstruction posterior  $p_\theta(\mathbf{A}|\mathbf{Z},\mathbf{X})$  to get our objective function.

# Importance of Reparameterization Trick

**Goal: Take Monte Carlo integrals to calculate the gradient of our objective.**

This goal might not be directly achievable because of problems such as...

- ▶ **No Analytic Solution**
- ▶ **Higher Variance**

**Summary:** The reparameterization trick can help in both departments.

## Reparameterization Trick: Case of No Analytic Solution

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p_{\theta}(z)}[f_{\theta}(z)] &= \nabla_{\theta} \left[ \int_z p_{\theta}(z) f_{\theta}(z) dz \right] \\ &= \int_z \nabla_{\theta} [p_{\theta}(z) f_{\theta}(z)] dz \\ &= \int_z f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz + \int_z p_{\theta}(z) \nabla_{\theta} f_{\theta}(z) dz \\ &= \underbrace{\int_z f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz}_{\text{What about this?}} + \mathbb{E}_{p_{\theta}(z)} [\nabla_{\theta} f_{\theta}(z)]\end{aligned}$$

Figure: Generic Gradient Calculation of Expected Value

Image Credit: <http://gregorygundersen.com/blog/2018/04/29/reparameterization/>

# Reparameterization Trick: Lowering Variance

In some cases ( $p_{\theta}(z) \sim N(\mu, \sigma^2)$  being one of them), our gradient can have an analytic solution, meaning that we can compute an expectation and we can estimate it using samples.

So... what's the problem?

Let's see the case I presented as an example!

## $q \sim N(\theta, 1)$ Example

Let's say that we were interested in the case where  $q_\theta \sim N(\theta, 1)$  and we are interested in the following optimization problem:

$$\theta = \operatorname{argmin}_\theta E_q[x^2]$$

Obviously we know the answer is  $\theta = 0$  because

$$\min_\theta E_q[x^2] = \min_\theta \left( (E_q[x])^2 + \operatorname{Var}(x) \right) = \min_\theta (\theta^2 + 1) \rightarrow \theta = 0$$

$\nabla_\theta E_q[x^2]$  is as follows

$$\nabla_\theta E_q[x^2] = \nabla_\theta \int q_\theta(x) x^2 dx = \int x^2 \nabla_\theta q_\theta(x) \frac{q_\theta(x)}{q_\theta(x)} dx = \int q_\theta(x) \nabla_\theta \log q_\theta(x) x^2 dx = E_q[x^2 \nabla_\theta \log q_\theta(x)]$$

if  $q_\theta(x) = N(\theta, 1)$ , this method gives

$$\nabla_\theta E_q[x^2] = E_q[x^2(x - \theta)]$$

Figure: Gradient Calculation of Normal Case

Image Credit:

# The Reparameterization

**Motivation:** The result on the previous slide shows we can sample from  $q$ ; so theoretically, everything is fine. Why would the reparameterization trick even be needed here?

Let's see what happens when we take the expectation with respect to a distribution that is independent of the parameter we took the gradient with respect to.

$$x = \theta + \epsilon, \quad \epsilon \sim N(0, 1)$$

$$E_q[x^2] = E_p[(\theta + \epsilon)^2]$$

can write the derivative of  $E_q[x^2]$  as follows

$$\nabla_{\theta} E_q[x^2] = \nabla_{\theta} E_p[(\theta + \epsilon)^2] = E_p[2(\theta + \epsilon)]$$

Figure: Reparameterization of Normal Case

# Reparameterization Simulation

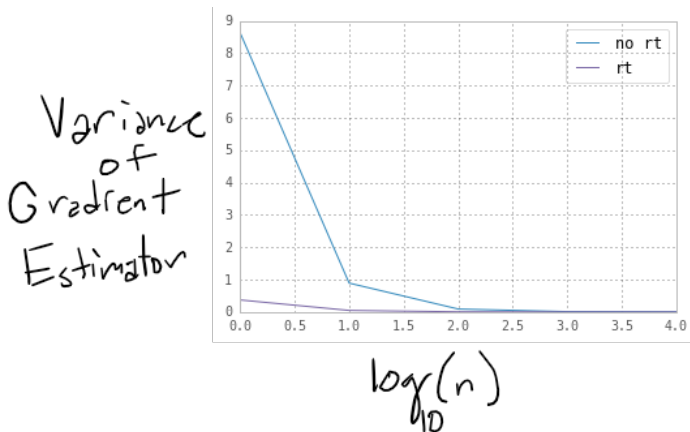


Figure: Variance of Gradient Estimator as a Function of Logarithmic Sample Size



# Formal Model

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$(\boldsymbol{\mu}, \boldsymbol{\sigma}) = \text{EncoderNN}_{\phi}(\mathbf{x})$$

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}))$$

$$(\boldsymbol{\mu}_l, \boldsymbol{\sigma}_l) = \text{DecoderNN}_{\theta}(\mathbf{z})$$

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_l, \text{diag}(\boldsymbol{\sigma}_l))$$

Image Credit: Ismael Mendoza

$$P(\mathbf{Z}|\mathbf{X}) \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

$$\boldsymbol{\mu}, \sigma = GNN_{\phi}(\mathbf{A}, \mathbf{X})$$

$$q_{\phi}(\mathbf{Z}|\mathbf{A}, \mathbf{X}) := \prod_{i=1}^n q_{\phi_i}(\mathbf{Z}_i|\mathbf{A}, \mathbf{X})$$

$$\hat{\mathbf{A}} = \frac{\mathbf{Z}\mathbf{Z}^T}{\|\mathbf{Z}\|^2} + \mathbf{1}\mathbf{1}^T$$

$$\mathbf{Z}^* = GNN_{\theta}(\hat{\mathbf{A}}, [\mathbf{Z}|\mathbf{X}])$$

$$p_{\theta}(\mathbf{A}|\mathbf{Z}, \mathbf{X}) = \prod_{i=1}^n p_{\theta}^{(i,j)}(\mathbf{A}_{ij}|\mathbf{Z}^*)$$

## Personal Tangent

Jeff and I are actually attempting to do something similar, but would use something like Graphite to learn a latent variable to generate  $\mathbf{X}$ , the feature matrix, instead of  $\mathbf{A}$ , the adjacency matrix. This tells us more about the nodal attributes of the graph while preserving graph relationships rather than trying to recreate the graph itself.

# Building Up Model Assumptions

**Property:** Under  $r(\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n | \mathbf{A}, \mathbf{X})$ , any  $\mathbf{Z}_i$  is independent of all other  $\mathbf{Z}_j$  when conditioned on  $\mathbf{A}, \mathbf{X}$  and all neighboring latent variables of node  $i$ .

**Takeaway:** Graphite usually expects latent distributions that follow this property so that the following optimization problem is solved and computationally inexpensive.

# Mean Field Approximation

A mean field approximation is an approximation that is used as follows:

$$r(\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n | \mathbf{A}, \mathbf{X}) \approx \prod_{i=1}^n q_{\phi_i}(\mathbf{Z}_i | \mathbf{A}, \mathbf{X})$$

# Variational Marginal Solution

**Objective:**

$$\min_{\phi_1, \dots, \phi_n} KL \left[ r(\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n | \mathbf{A}, \mathbf{X}) || \prod_{i=1}^n q_{\phi_i}(\mathbf{Z}_i | \mathbf{A}, \mathbf{X}) \right]$$

**Solution:**  $q_{\phi_i}(\mathbf{Z}_i | \mathbf{A}, \mathbf{X}) = \mathcal{O}_{\mathcal{G}}^{MF}(\mathbf{Z}_i, \{q_{\phi_j}\}_{j \in \mathcal{N}_{(i)}})$  where  $\mathcal{N}_{(i)}$  is the neighbors of node  $i$

This solution is provided by Wainwright et al. 2008 and is given without proof.

# Connection Between WL and Kernel Embeddings with GNN (Eq 17 & Theorem 2)

$$q_{\phi_i}^{(l)}(\mathbf{Z}_i | \mathbf{A}, \mathbf{X}) = \mathcal{O}_{\mathcal{G}}^{MF}(\mathbf{Z}_i, \{q_{\phi_j}^{(l-1)}\}_{j \in \mathcal{N}(i)})$$

$$\mu_i^{(l)} = \mathcal{O}_{\mathcal{G}}^{MF}(\{\mu_i^{(l)}\}_{j \in \mathcal{N}(i)})$$

**Theorem 2.** Let  $\mathcal{G}$  be any undirected latent variable model such that the conditional distribution  $r(\mathbf{Z}_1, \dots, \mathbf{Z}_n | \mathbf{A}, \mathbf{X})$  expressed by the model satisfies Property 1.

Then there exists a choice of  $\eta_l$ ,  $\mathcal{F}_l$ ,  $\mathbf{W}_l$ , and  $\mathbf{B}_l$  such that for all  $\{\mu_i^{(l-1)}\}_{i=1}^n$ , the GNN propagation rule in Eq. (2) is computationally equivalent to updating  $\{\mu_i^{(l-1)}\}_{i=1}^n$  via a first order approximation of Eq. (17).

*Proof.* See Appendix A.



Figure

Image Credit: *Graphite: Iterative Generative Modeling of Graphs*

## How GNNs Accomplish This (Last Sentence of Sec. 5)

This is why GNNs in a VAE setting are so powerful. They can learn the parameters needed to make the computation incredibly simple while still doing the impressive learning that a VAE does.

**Note:** I am just relaying the result Theorem 2 back to you. I'm not entirely sure how this result is practical in the paper and have left it as a discussion question.

## Experimental Results: Reconstruction

Graphite competes against (Variational) Graph Autoencoder. The task is to reconstruct the original graph (autoencoder case) or create a good latent variable model (variational autoencoder case).

Table 1. Mean reconstruction errors and negative log-likelihood estimates (in nats) for autoencoders and variational autoencoders respectively on test instances from six different generative families. Lower is better.

	Erdo-Renyi	Ego	Regular	Geometric	Power Law	Barabasi-Albert
GAE	221.79 $\pm$ 7.58	197.3 $\pm$ 1.99	198.5 $\pm$ 4.78	514.26 $\pm$ 41.58	519.44 $\pm$ 36.30	236.29 $\pm$ 15.13
Graphite-AE	<b>195.56</b> $\pm$ 1.49	<b>182.79</b> $\pm$ 1.45	<b>191.41</b> $\pm$ 1.99	<b>181.14</b> $\pm$ 4.48	<b>201.22</b> $\pm$ 2.42	<b>192.38</b> $\pm$ 1.61
VGAE	273.82 $\pm$ 0.07	273.76 $\pm$ 0.06	275.29 $\pm$ 0.08	274.09 $\pm$ 0.06	278.86 $\pm$ 0.12	274.4 $\pm$ 0.08
Graphite-VAE	<b>270.22</b> $\pm$ 0.15	<b>270.70</b> $\pm$ 0.32	<b>266.54</b> $\pm$ 0.12	<b>269.71</b> $\pm$ 0.08	<b>263.92</b> $\pm$ 0.14	<b>268.73</b> $\pm$ 0.09

**Figure:** Results of (V)GAE v. Graphite in Reconstruction and Density Estimation

Image Credit: *Graphite: Iterative Generative Modeling of Graphs*

Why the overwhelmingly better performance? GAE/VGAE doesn't use learnable parameters for its graph reconstruction where Graphite does. Graphite learns  $\theta$  for  $\mathbf{Z}^*$  iteratively until the final graph reconstruction.



# Experimental Results: Link Prediction

**Goal:** Predict whether two nodes share an edge.

**Setup:**

1. Given a complete graph, take a connected subgraph that has 85% of the edges of the original graph. This removes some true edges from the data.
2. Take the 85% subgraph for training, 5% of the removed edges for validation, and 10% of the removed edges for testing.
3. Balance the validation and testing sets by injecting some fake edges. This allows for our model to have the option to make some false positive classifications.
4. Run Graphite!

# Experimental Results: Link Prediction (cont.)

Table 3. Area Under the ROC Curve (AUC) for link prediction (\* denotes dataset with features). Higher is better.

	Cora	Citeseer	Pubmed	Cora*	Citeseer*	Pubmed*
SC	89.9 $\pm$ 0.20	91.5 $\pm$ 0.17	<b>94.9 <math>\pm</math> 0.04</b>	-	-	-
DeepWalk	85.0 $\pm$ 0.17	88.6 $\pm$ 0.15	91.5 $\pm$ 0.04	-	-	-
node2vec	85.6 $\pm$ 0.15	89.4 $\pm$ 0.14	91.9 $\pm$ 0.04	-	-	-
GAE	90.2 $\pm$ 0.16	92.0 $\pm$ 0.14	92.5 $\pm$ 0.06	93.9 $\pm$ 0.11	94.9 $\pm$ 0.13	96.8 $\pm$ 0.04
VGAE	90.1 $\pm$ 0.15	92.0 $\pm$ 0.17	92.3 $\pm$ 0.06	94.1 $\pm$ 0.11	96.7 $\pm$ 0.08	95.5 $\pm$ 0.13
Graphite-AE	91.0 $\pm$ 0.15	92.6 $\pm$ 0.16	94.5 $\pm$ 0.05	94.2 $\pm$ 0.13	96.2 $\pm$ 0.10	<b>97.8 <math>\pm</math> 0.03</b>
Graphite-VAE	<b>91.5 <math>\pm</math> 0.15</b>	<b>93.5 <math>\pm</math> 0.13</b>	94.6 $\pm$ 0.04	<b>94.7 <math>\pm</math> 0.11</b>	<b>97.3 <math>\pm</math> 0.06</b>	97.4 $\pm$ 0.04

Table 4. Average Precision (AP) scores for link prediction (\* denotes dataset with features). Higher is better.

	Cora	Citeseer	Pubmed	Cora*	Citeseer*	Pubmed*
SC	92.8 $\pm$ 0.12	94.4 $\pm$ 0.11	<b>96.0 <math>\pm</math> 0.03</b>	-	-	-
DeepWalk	86.6 $\pm$ 0.17	90.3 $\pm$ 0.12	91.9 $\pm$ 0.05	-	-	-
node2vec	87.5 $\pm$ 0.14	91.3 $\pm$ 0.13	92.3 $\pm$ 0.05	-	-	-
GAE	92.4 $\pm$ 0.12	94.0 $\pm$ 0.12	94.3 $\pm$ 0.5	94.3 $\pm$ 0.12	94.8 $\pm$ 0.15	96.8 $\pm$ 0.04
VGAE	92.3 $\pm$ 0.12	94.2 $\pm$ 0.12	94.2 $\pm$ 0.04	94.6 $\pm$ 0.11	97.0 $\pm$ 0.08	95.5 $\pm$ 0.12
Graphite-AE	92.8 $\pm$ 0.13	94.1 $\pm$ 0.14	95.7 $\pm$ 0.06	94.5 $\pm$ 0.14	96.1 $\pm$ 0.12	<b>97.7 <math>\pm</math> 0.03</b>
Graphite-VAE	<b>93.2 <math>\pm</math> 0.13</b>	<b>95.0 <math>\pm</math> 0.10</b>	<b>96.0 <math>\pm</math> 0.03</b>	<b>94.9 <math>\pm</math> 0.13</b>	<b>97.4 <math>\pm</math> 0.06</b>	97.4 $\pm$ 0.04

Figure: Results of Link Prediction

Image Credit: *Graphite: Iterative Generative Modeling of Graphs*

## Experiment Results: Semi-Supervised Node Learning

**Goal:** Given labels of some nodes in a graph, try to classify some of the unlabeled nodes.

Table 5. Classification accuracies (\* denotes dataset with features).  
Baseline numbers from Kipf & Welling (2017).

	Cora*	Citeseer*	Pubmed*
SemiEmb	59.0	59.6	71.1
DeepWalk	67.2	43.2	65.3
ICA	75.1	69.1	73.9
Planetoid	75.7	64.7	77.2
GCN	81.5	70.3	79.0
Graphite	<b>82.1 ± 0.06</b>	<b>71.0 ± 0.07</b>	<b>79.3 ± 0.03</b>

Figure: Results of Node Classification

Image Credit: *Graphite: Iterative Generative Modeling of Graphs*

**Important Takeaway:** Graphite is incredibly strong for this task because it has a default behavior that is transductive. This means that Graphite can use testing nodes without their labels. This is because even if Graphite doesn't have access to a node label, it has access to its locality (node neighbors). Therefore,  $f(\mathbf{A})$  still contains a lot of information for unlabeled nodes.