

Table 3: Statistics of the datasets used in graph feature imputation

Dataset	Nodes	Train/Test Features	Feature Density	Edges	Edge Density
Douban	3,000	123,202/13,689	1.52%	2,690	0.03%
Ciao	7,317	39,279/16,892	0.77%	111,781	0.21%
Cora	2,708	88,590/9,842	2.54%	5,429	0.15%
Citeseer	3,327	189,298/21,032	1.70%	4,732	0.08%
Amaphoto	7,650	35,640/3,958	0.69%	143,663	0.49%
Amacomp	13,752	66,150/7,350	0.69%	287,209	0.30%

Table 4: Statistics of the datasets used in graph structure generation

Datasets	MUTAG	PTC-MR	ZINC
No.Graphs	188	344	250k
Avg.Nodes	17.9	14.3	23.1
Avg.Edges	19.8	14.7	23.9

464 A Proof of Proposition 1

465 *Proof.* Based on the independent white noise model (7) and (8), we have

$$\begin{aligned}
x' &= U \text{diag}(g_c^{-1}(\lambda)) U^\top \sigma^{-1}(\hat{h}) \\
&= U \text{diag}(g_c^{-1}(\lambda)) U^\top \sigma^{-1}(\sigma(U \text{diag}(g_c(\lambda)) U^\top x) + \epsilon) \\
&= x + U \text{diag}(g_c^{-1}(\lambda)) U^\top \sigma^{-1}(\epsilon).
\end{aligned}$$

466 Hence,

$$(x' - x) = U \text{diag}(g_c^{-1}(\lambda)) U^\top \sigma^{-1}(\epsilon). \quad (14)$$

467 In general, we cannot find an exact formula for the variance of a nonlinear function. Here, we
468 consider a more relevant case, that is, a random variable ϵ is fairly closed to the mean (i.e., the
469 standard error is sufficiently small), then we can approximate the nonlinear transformation $\sigma^{-1}(\epsilon)$
470 with a first-order Taylor expansion. Under the condition $\sigma^{-1}(\epsilon) \approx \sigma^{-1}(\mathbf{0}) + \nabla \sigma^{-1}(\mathbf{0})^T \epsilon$ where
471 $\sigma^{-1}(\cdot)$ is a element-wise function. Plugging into (14), we have

$$\begin{aligned}
(x' - x) &\approx U \text{diag}(g_c^{-1}(\lambda)) U^\top (\sigma^{-1}(\mathbf{0}) + \nabla \sigma^{-1}(\mathbf{0})^T \epsilon) \\
x' &\approx U \text{diag}(g_c^{-1}(\lambda)) U^\top \nabla \sigma^{-1}(\mathbf{0})^T \epsilon + U \text{diag}(g_c^{-1}(\lambda)) U^\top \sigma^{-1}(\mathbf{0}) + x
\end{aligned} \quad (15)$$

472 Hence, we can observe that the variance of x' is almost identical to the variance of
473 $U \text{diag}(g_c^{-1}(\lambda)) U^\top \nabla \sigma^{-1}(\mathbf{0})^T \epsilon$ in an asymptotic sense.

$$\begin{aligned}
\text{Var}(x') &\approx (\partial \sigma^{-1})^2(0) \text{Var}(U \text{diag}(g_c^{-1}(\lambda)) U^\top \epsilon) \\
&\geq \text{Var}(U \text{diag}(g_c^{-1}(\lambda)) U^\top \epsilon).
\end{aligned} \quad (16)$$

474 The last inequality holds as our assumption on the inverse activation function. We will elaborate more
475 details and explain why it is a reasonable condition.

476 **Fact 2.** If z is a Gaussian random variable, i.e., $z \sim \mathcal{N}(\mu, \Sigma)$, then $Az \sim \mathcal{N}(A\mu, A\Sigma A^T)$.

$$\begin{aligned}
\Delta &= U \text{diag}(g_c^{-1}(\lambda)) U^\top \epsilon \\
&= U \text{diag}(g_c^{-1}(\lambda)) \xi \quad (\text{denote } \xi = U^\top \epsilon).
\end{aligned}$$

477 Here, the distribution of ξ is same as ϵ 's due to the orthogonal matrix U and the rotation invariance of
478 the Gaussian distribution. We refer the reader to Fact 2 for details. Hence, note that $y = \text{diag}(g_c^{-1}(\lambda)) \xi$,
479 we have

$$y \sim \mathcal{N}(0, \sigma^2 \text{diag}([(1 - \lambda_1)^{-2}, \dots, (1 - \lambda_N)^{-2}])).$$

480

$$\text{Var}(x') = \text{Var}(\Delta) = \sigma^2 \sum_{i=1}^N \frac{1}{(1 - \lambda_i)^2}.$$

481 Since the eigenvalue of the normalized Laplacian matrix is bounded, that is, $\lambda_i \in [0, 2], \forall i \in [N]$, we
 482 have $1 - \lambda_i \in [-1, 1]$ and $\frac{1}{(1 - \lambda_i)^2} \in [1, +\infty)$. It was worthwhile mentioning that the equality holds
 483 only when $\lambda_i = 0$ or $2, \forall i \in [N]$. Therefore,

$$\text{Var}(x') = \sigma^2 \sum_{i=1}^N \frac{1}{(1 - \lambda_i)^2} \geq N\sigma^2 = \text{Var}(\epsilon).$$

484 Regarding the condition $\partial\sigma^{-1}(0) \geq 1$, we can consider three representative examples to corroborate
 485 our assumption, i.e.,

486 • Sigmoid function; $\sigma(x) = \frac{1}{1+e^{-x}}$ and its inverse operator $\sigma^{-1}(x) = -\ln\left(\frac{1}{x} - 1\right)$. Thus,
 487 we have

$$\partial\sigma^{-1}(x) = \frac{1}{x - x^2},$$

488 and $(\partial\sigma^{-1}(x))^2 \rightarrow +\infty$ when x goes to zero.

489 • Hyperbolic tangent (Tanh) activation function; $\sigma(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$. The inverse hyperbolic
 490 tangent is given as

$$\sigma^{-1}(x) = \frac{1}{2}[\ln(1+x) - \ln(1-x)].$$

491 Computing its gradient, we have

$$\partial\sigma^{-1}(x) = \frac{1}{1 - x^2} \geq 1.$$

• Leaky ReLU;

$$\sigma(x) = \begin{cases} x & \text{if } x > 0 \\ \frac{x}{\alpha} & \text{if } x \leq 0, \end{cases} \Rightarrow \sigma^{-1}(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x \leq 0, \end{cases}$$

492 where α is a constant belonging to $(1, +\infty)$. Therefore, the subgradient of σ^{-1} is larger or
 493 equal to 1.

494

□

495 **Remark 3.** *The last equality is extremely hard to hold in practice and indeed to corroborate our*
 496 *empirical finding that the inverse operator will amplify the noise.*

497 This concludes the proof of proposition 1. Next, we show proposition 1 still holds if GCN model is
 498 parameterized.

499 **Parameterized Model** We consider

$$\hat{h} = \sigma(WU \text{diag}(g_c(\lambda))U^\top x) + \epsilon, \quad (17)$$

500 where ϵ is the white Gaussian noise (i.e., $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_N)$). Under the condition that W is orthogonal,
 501 the inverse operator admits

$$x' = U \text{diag}(g_c^{-1}(\lambda))U^\top W^\top \sigma^{-1}(\hat{h}). \quad (18)$$

502 Therefore,

$$\begin{aligned} x' &= U \text{diag}(g_c^{-1}(\lambda))U^\top W^\top \sigma^{-1}(\hat{h}) \\ &= U \text{diag}(g_c^{-1}(\lambda))U^\top W^\top \sigma^{-1}(\sigma(WU \text{diag}(g_c(\lambda))U^\top x) + \epsilon) \\ &= x + U \text{diag}(g_c^{-1}(\lambda))U^\top W^\top \sigma^{-1}(\epsilon). \end{aligned}$$

503 It was worth noting that the product of orthogonal matrices is orthogonal. By the similar argument in
 504 Proposition 1, we can also conclude that the proposed inverse operation may introduce undesirable
 505 noise into the output graph.

506 B Comparisons with Additional Baselines

507 We additionally provide the comparison results of our method with three discrete matrix completion
508 models: sRGCNN [29], GC-MC [4], and GraphRec [12]. We adapt the same setting as in Section
509 4.1 and the results are shown in Table 5. Ours yields the smallest imputation error on all datasets
510 compared with the three matrix completion models. Notice that the three matrix completion models
511 are specially designed for discrete matrix completion, where ours is applicable to both continuous
512 and discrete feature values.

513 C Graph Autoencoder Framework

514 C.1 Encoder

515 Our encoder consists of two layers of Graph Convolutional Networks (GCN) [21], to produce
516 smoothed representations of the input graphs.

Convolution

$$H = \text{GCN}(A, X), \quad (19)$$

517 where $H \in \mathbb{R}^{N \times v}$ denotes smoothed node representations.

518 C.2 Decoder

519 Our decoder consists of one layer of GDN, to produce fine-grained graphs.

Deconvolution

$$X' = \text{GDN}(A, H), \quad (20)$$

520 where $X' \in \mathbb{R}^{N \times d}$ denotes the recovered graph features.

521 C.3 The loss function

522 The reconstruction loss is a feature reconstruction loss.

$$\mathcal{L} = f(X, X'), \quad (21)$$

523 where $f(\cdot, \cdot)$ denotes a differential distance metric, e.g., $f(\cdot, \cdot) = \text{MSE}(\cdot, \cdot)$, and $\text{MSE}(\cdot, \cdot)$ represents
524 Mean Squared Error.

525 D Detailed Model Configuration

526 All the experiments are done on a single machine with 32 cores and 300G memories. We don't use
527 GPU as it involves a lot of sparse matrix multiplication.

528 For graph feature imputation, We use the graph autoencoder framework specified in Appendix C.
529 The recovered graph features X' is used to infer the potential missing features. We train our model
530 using Adam optimizer with a learning rate of $\{0.005, 0.002\}$. The output dimension of the first layer
531 is $\{256, 512\}$. The output dimension of the second layer is 128. The number of epochs is chosen
532 from $\{100, 200\}$. We use full-batch size. We stack the output of the first GCN layer and the second
533 GCN layer in our encoders, and use left normalization to preprocess adjacency matrix. We don't use
534 feature Dropout as it does not improve the performance. We use DropEdge [33]. We use grid search
535 to choose the hyper-parameters. The selected parameters are reported in Table 6.

536 For graph generation, We strictly follow the architectures of VGAE [20] and Graphite [15] and add
537 one feature reconstruction loss. The overall reconstruction loss is a sum of structure reconstruction
538 term and feature reconstruction term,

$$\log p(G|Z) = \log p(A|Z) + \log p(X|Z), \quad (22)$$

539 where $p(A|Z)$ is the default structure reconstruction term such as a simple inner product of latent
540 variables in VGAE, i.e., $\prod_i \prod_j p(A_{ij}|z_i, z_j)$. $\log p(X|Z)$ is computed with GDN and defined as

Table 5: RMSE test comparison with matrix completion methods on graph feature imputation

Datasets	Ciao	Douban	Cora	Citeseer	Amaphoto	Amacomp
sRGCNN [29]	1.183	0.801	0.550	0.551	0.519	0.509
GC-MC [4]	1.061	0.734	0.434	0.420	0.402	0.405
GraphRec [12]	1.062	0.754	0.437	0.425	0.417	0.417
OURS	1.011	0.734	0.415	0.399	0.391	0.393

Table 6: Summary of chosen hyper-parameters on graph feature imputation

Datasets	Ciao	Douban	Cora	Citeseer	Amaphoto	Amacomp
Dim of the first layer	256	256	512	256	256	256
Dim of the second layer	128	128	64	128	128	128
Learning rate	0.005	0.005	0.002	0.005	0.005	0.005
Epochs	200	200	200	200	100	100
DropEdge rate	1	0.5	0.5	0.5	1	1

541 in Eq. 21. As both VGAE and Graphite can only deal with homogeneous edge attribute, we don't
 542 differentiate bond types in ZINC dataset, i.e., we only evaluate if there is a bond between two given
 543 atoms. The number of epochs, initial learning rate and hyperparameters are all set to be the same
 544 with VGAE [20] or Graphite [15], for a fair comparison. The effect of orders of Maclaurin Series in
 545 GDN are reported in Table 7.

546 E Derivative of Inverse Operator

$$g_c^{-1}(\lambda_i) = \frac{1}{1 - \lambda_i} = \sum_{n=0}^{\infty} \frac{\left(\frac{1}{1-\lambda_i}\right)_{\lambda_i=0}^{(n)}}{n!} \lambda_i^n = \sum_{n=0}^{\infty} \frac{(-1)^n n! (-1)^n}{n!} \lambda_i^n = \sum_{n=0}^{\infty} \lambda_i^n$$

547 where (n) denotes the n -th order derivative. Thus, Eq. 9 can be obtained as:

$$U \text{diag}\left(\sum_{n=0}^{\infty} \lambda_1^n, \dots, \sum_{n=0}^{\infty} \lambda_N^n\right) U^T = U \sum_{n=0}^{\infty} \Lambda^n U^T \quad (23)$$

Table 7: Summary of the effect of orders of Maclaurin Series in GDN on the performance of graph structure generation

Datasets		MUTAG			PTC-MR			ZINC		
-		$\log p(A Z)$	AUC	AP	$\log p(A Z)$	AUC	AP	$\log p(A Z)$	AUC	AP
VGAE + GDN	First-order	-1.002	0.738	0.464	-1.351	0.493	0.386	-0.999	0.474	0.239
	Second-order	-1.008	0.766	0.502	-1.356	0.559	0.443	-1.001	0.611	0.345
	Third-order	-1.026	0.823	0.611	-1.351	0.760	0.602	-1.006	0.858	0.611
Graphite + GDN	First-order	-1.003	0.741	0.470	-1.355	0.529	0.398	-1.001	0.478	0.241
	Second-order	-1.014	0.774	0.501	-1.358	0.574	0.450	-1.003	0.632	0.379
	Third-order	-1.024	0.818	0.608	-1.347	0.773	0.613	-0.998	0.838	0.567