
On the Unreasonable Effectiveness of Feature Propagation in Learning on Graphs with Missing Node Features

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 While Graph Neural Networks (GNNs) have recently become the *de facto* standard
2 for modeling relational data, they impose a strong assumption on the availability of
3 the node or edge features of the graph. In many real-world applications, however,
4 features are only partially available; for example, in social networks, age and
5 gender are available only for a small subset of users. We present a general approach
6 for handling missing features in graph machine learning applications that is based
7 on minimization of the Dirichlet energy and leads to a diffusion-type differential
8 equation on the graph. The discretization of this equation produces a simple, fast
9 and scalable algorithm which we call Feature Propagation. We experimentally show
10 that the proposed approach outperforms previous methods on seven common node-
11 classification benchmarks and can withstand surprisingly high rates of missing
12 features: on average we observe only around 4% relative accuracy drop when 99%
13 of the features are missing. Moreover, it takes only 10 seconds to run on a graph
14 with $\sim 2.5\text{M}$ nodes and $\sim 123\text{M}$ edges on a single GPU.

15 1 Introduction

16 Graph Neural Networks (GNNs) [6, 19, 21, 30, 40, 47] have been successful on a broad range
17 of problems and in a variety of fields [13, 14, 17, 38, 43, 54, 60]. GNNs typically operate by a
18 message-passing mechanism [3, 20], where at each layer, nodes send their feature representations
19 (“messages”) to their neighbors. The feature representation of each node is initialized to their original
20 features, and is updated by repeatedly aggregating incoming messages from neighbors. Being able to
21 combine the topological information with feature information is what distinguishes GNNs from other
22 purely topological learning approaches such as random walks [22, 36] or label propagation [58], and
23 arguably what leads to their success.

24 GNN models typically assume a fully observed feature matrix, where rows represent nodes and
25 columns feature channels. However, in real-world scenarios, each feature is often only observed for a
26 subset of the nodes. For example, demographic information can be available for only a small subset
27 of social network users, while content features are generally only present for the most active users. In
28 a co-purchase network, not all products may have a full description associated with them. With the
29 rising awareness around digital privacy, data is increasingly available only upon explicit user consent.
30 In all the above cases, the feature matrix contains missing values and most existing GNN models
31 cannot be directly applied.

32 While classic imputation methods [29, 32, 55] can be used to fill the missing values of the feature
33 matrix, they are unaware of the underlying graph structure. Graph Signal Processing, a field attempting
34 to generalize classical Fourier analysis to graphs, offers several methods that reconstruct signals on

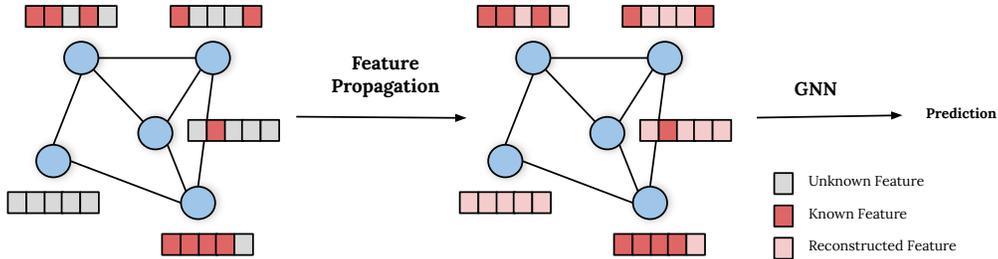


Figure 1: A diagram illustrating our Feature Propagation framework. On the left, a graph with missing node features. In the initial reconstruction step, Feature Propagation reconstructs the missing features by iteratively diffusing the known features in the graph. Subsequently, the graph and the reconstructed node features are fed into a downstream GNN model, which then produces a prediction.

35 graphs [34]. However, they do not scale beyond graphs with a few thousand nodes, making them
 36 infeasible for practical applications. More recently, SAT [10], GCNMF [44] and PaGNN [26] have
 37 been proposed to adapt GNNs to the case of missing features. However, they are not evaluated at high
 38 missing features rates ($> 90\%$), which occur in many real-world scenarios, and where we find them
 39 to suffer. Moreover, they are unable to scale to graphs with more than a few hundred thousand nodes.
 40 At the time of writing, PaGNN is the state-of-the-art method for node classification with missing
 41 features.

42 **Contributions** We present a general approach for handling missing node features in graph machine
 43 learning tasks. The framework consists of an initial diffusion-based feature reconstruction step
 44 followed by a downstream GNN. The reconstruction step is based on Dirichlet energy minimization,
 45 which leads to a diffusion-type differential equation on the graph. Discretization of this differential
 46 equation leads to a very simple, fast, and scalable iterative algorithm which we call Feature
 47 Propagation (FP). FP outperforms state-of-the-art methods on six standard node-classification
 48 benchmarks and presents the following advantages:

- 49 • **Theoretically Motivated:** FP emerges naturally as the gradient flow minimizing the Dirichlet
 50 energy and can be interpreted as a diffusion equation on the graph with known features used as
 51 boundary conditions. This contributes to the promising direction of building continuous-time
 52 models on graphs.
- 53 • **Robust to high rates of missing features:** FP can withstand surprisingly high rates of missing
 54 features. In our experiment, we observe on average around 4% relative accuracy drop when up to
 55 99% of the features are missing. In comparison, GCNMF and PaGNN have an average drop of
 56 53.33% and 21.25% respectively. This finding has important implications especially in scenarios
 57 where the cost of sampling (observing features on nodes) is high or sampling is not possible
 58 altogether.
- 59 • **Generic:** FP can be combined with any GNN model to solve the downstream task; in contrast,
 60 GCNMF and PaGNN are specific GCN-type models.
- 61 • **Fast and Scalable:** FP takes only around 10 seconds for the reconstruction step on OGBN-
 62 Products (a graph with $\sim 2.5\text{M}$ nodes and $\sim 123\text{M}$ edges) on a single GPU. GCNMF and PaGNN
 63 run out-of-memory on this dataset.

65 2 Preliminaries

66 Let $G = (V, E)$ be an undirected graph with $n \times n$ adjacency matrix \mathbf{A} and a node feature vector¹
 67 $\mathbf{x} \in \mathbb{R}^n$. The *graph Laplacian* is an $n \times n$ positive semi-definite matrix $\mathbf{\Delta} = \mathbf{I} - \tilde{\mathbf{A}}$, where
 68 $\tilde{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$ is the normalized adjacency matrix and $\mathbf{D} = \text{diag}(\sum_j a_{1j}, \dots, \sum_j a_{nj})$ is the
 69 diagonal degree matrix.

¹For convenience, we assume scalar node features. Our derivations apply straightforwardly to the case of d -dimensional features represented as an $n \times d$ matrix \mathbf{X} .

Denote by $V_k \subseteq V$ the set of nodes on which the features are *known*, and by $V_u = V_k^c = V \setminus V_k$ the *unknown* ones. We further assume the ordering of the nodes such that we can write

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_u \end{bmatrix} \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_{kk} & \mathbf{A}_{ku} \\ \mathbf{A}_{uk} & \mathbf{A}_{uu} \end{bmatrix} \quad \Delta = \begin{bmatrix} \Delta_{kk} & \Delta_{ku} \\ \Delta_{uk} & \Delta_{uu} \end{bmatrix}.$$

70 Because the graph is undirected, \mathbf{A} is symmetric and thus $\mathbf{A}_{ku}^\top = \mathbf{A}_{uk}$ and $\Delta_{ku}^\top = \Delta_{uk}$. We will
71 tacitly assume this in the following discussion.

72 **Graph feature interpolation** is the problem of reconstructing the unknown features \mathbf{x}_u given the
73 graph structure G and the known features \mathbf{x}_k . The interpolation task requires some prior on the
74 behavior of the features of the graph, which can be expressed in the form of an energy function
75 $\ell(\mathbf{x}, G)$. The most common assumption is feature *homophily* (i.e., that the features of every node are
76 similar to those of the neighbours), quantified using a criterion of *smoothness* such as the Dirichlet
77 energy. Since in many cases the behavior of the features is not known, the energy can possibly be
78 learned from the data.

79 **Learning on a graph with missing features** is a transductive learning problem (typically node-
80 wise classification or regression using some GNN architecture) where the structure of the graph
81 G is known while the labels and node features are only partially known on the subsets V_l and V_k
82 of nodes, respectively (that might be different and even disjoint). Specifically, we try to learn a
83 function $\mathbf{f}(\mathbf{x}_k, G)$ such that $f_i \approx y_i$ for $i \in V_l$. Learning with missing features can be done by a
84 pre-processing step of graph signal interpolation (reconstructing an estimate $\tilde{\mathbf{x}}$ of the full feature
85 vector \mathbf{x} from \mathbf{x}_k) independent of the learning task, followed by the learning task of $\mathbf{f}(\tilde{\mathbf{x}}, G)$ on the
86 inferred fully-featured graph. In some settings, we are not interested in recovering the features *per se*,
87 but rather ensuring that the output of the *function* \mathbf{f} on these features is correct – arguably a more
88 ‘forgiving’ setting.

89 3 Feature Propagation

90 We assume to be given \mathbf{x}_k and attempt to find the missing node features \mathbf{x}_u by means of interpolation
91 that minimizes some energy $\ell(\mathbf{x}, G)$. In particular, we consider the *Dirichlet energy* $\ell(\mathbf{x}, G) =$
92 $\frac{1}{2} \mathbf{x}^\top \Delta \mathbf{x} = \frac{1}{2} \sum_{ij} \tilde{a}_{ij} (x_i - x_j)^2$, where \tilde{a}_{ij} are the individual entries of the normalized adjacency
93 $\tilde{\mathbf{A}}$. The Dirichlet energy is widely used as a smoothness criterion for functions defined on the nodes
94 of the graph and thus promotes homophily. Functions minimizing the Dirichlet energy are called
95 *harmonic*; without boundary conditions, it is minimized by a constant function.

96 While the Dirichlet energy is convex and it is possible to derive its minimizer in a closed-form, as
97 shown in Appendix A.2, its computational complexity makes it unfeasible for graphs with many
98 nodes with missing features. Instead, we consider the associated *gradient flow* $\dot{\mathbf{x}}(t) = -\nabla \ell(\mathbf{x}(t))$
99 as a differential equation with boundary condition $\mathbf{x}_k(t) = \mathbf{x}_k$ whose solution at the missing nodes,
100 $\mathbf{x}_u = \lim_{t \rightarrow \infty} \mathbf{x}_u(t)$, provides the desired interpolation.

Gradient flow. For the Dirichlet energy, $\nabla_{\mathbf{x}} \ell = \Delta \mathbf{x}$ and the gradient flow takes the form of the
standard isotropic heat diffusion equation on the graph,

$$\dot{\mathbf{x}}(t) = -\Delta \mathbf{x}(t) \quad (\text{IC}) \quad \mathbf{x}(0) = \begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_u(0) \end{bmatrix} \quad (\text{BC}) \quad \mathbf{x}_k(t) = \mathbf{x}_k$$

101 where IC and BC stand for initial conditions and boundary conditions respectively. By incorporating
102 the boundary conditions, we can compactly express the diffusion equation as

$$\begin{bmatrix} \dot{\mathbf{x}}_k(t) \\ \dot{\mathbf{x}}_u(t) \end{bmatrix} = - \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \Delta_{uk} & \Delta_{uu} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_u(t) \end{bmatrix} = - \begin{bmatrix} \mathbf{0} \\ \Delta_{uk} \mathbf{x}_k + \Delta_{uu} \mathbf{x}_u(t) \end{bmatrix}. \quad (1)$$

103 As expected, the gradient flow of the observed features is $\mathbf{0}$, given that they do not change during the
104 diffusion.

105 The evolution of the missing features can be regarded as a heat diffusion equation with a constant
106 heat source $\Delta_{uk} \mathbf{x}_k$ coming from the boundary (known) nodes. Since the graph Laplacian matrix is

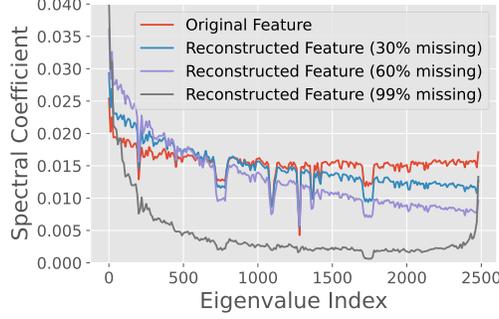


Figure 2: Graph Fourier transform magnitudes of the original Cora features (red) and those reconstructed by FP for varying rates of missing rates (we take the average over feature channels). Since FP minimizes the Dirichlet energy, it can be interpreted as a low-pass filter, which is stronger for a higher rate of missing features.

107 positive semi-definite, the Dirichlet energy ℓ is convex. Its global minimizer is given by the solution
 108 to the closed-form equation $\nabla_{\mathbf{x}_u} \ell = \mathbf{0}$ and by rearranging the final $|V_u|$ rows of Equation 1 we get
 109 the solution $\mathbf{x}_u = -\Delta_{uu}^{-1} \Delta_{ku}^\top \mathbf{x}_k$. This solution always exists as Δ_{uu} is non-singular, by virtue of
 110 the following:

111 **Proposition 3.1** (The sub-Laplacian matrix of an undirected connected graph is invertible). *Take*
 112 *any undirected, connected graph with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$, and its Laplacian $\Delta =$*
 113 *$\mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$, with \mathbf{D} being the degree matrix of \mathbf{A} . Then, for any principle sub-matrix*
 114 *$\mathbf{L}_u \in \mathbb{R}^{b \times b}$ of the Laplacian, where $1 \leq b < n$, \mathbf{L}_u is invertible.*

115 *Proof:* See Appendix A.2. Also, while the proposition assumes that the graph is connected, our
 116 analysis and method generalize straightforwardly in the case of a disconnected graph as we can
 117 simply apply Feature Propagation to each connected component independently.

118 However, solving a system of linear equations is computationally expensive (incurring $\mathcal{O}(|V_u|^3)$
 119 complexity for matrix inversion) and thus intractable for anything but only small graphs.

120 **Iterative scheme.** As an alternative, we can discretize the diffusion equation (1) and solve it by an
 121 iterative numerical scheme. Approximating the temporal derivative as forward difference with the
 122 time variable t discretized using a fixed step ($t = hk$ for step size $h > 0$ and $k = 1, 2, \dots$), we obtain
 123 the *explicit Euler scheme*:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - h \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \Delta_{uk} & \Delta_{uu} \end{bmatrix} \mathbf{x}^{(k)} = \left(\mathbf{I} - \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ h\Delta_{uk} & h\Delta_{uu} \end{bmatrix} \right) \mathbf{x}^{(k)} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -h\Delta_{uk} & \mathbf{I} - h\Delta_{uu} \end{bmatrix} \mathbf{x}^{(k)}$$

For the special case of $h = 1$, we can use the following observation

$$\tilde{\mathbf{A}} = \mathbf{I} - \Delta = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} - \begin{bmatrix} \Delta_{kk} & \Delta_{ku} \\ \Delta_{uk} & \Delta_{uu} \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \Delta_{kk} & -\Delta_{ku} \\ -\Delta_{uk} & \mathbf{I} - \Delta_{uu} \end{bmatrix},$$

124 to write the iteration formula as

$$\mathbf{x}^{(k+1)} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \tilde{\mathbf{A}}_{uk} & \tilde{\mathbf{A}}_{uu} \end{bmatrix} \mathbf{x}^{(k)}. \quad (2)$$

125 The Euler scheme is the gradient descent of the Dirichlet energy. Thus, applying the scheme decreases
 126 the Dirichlet energy and results in the features becoming increasingly smooth. Iteration (2) can be
 127 interpreted as successive low-pass filtering. Figure 2 depicts the magnitude of the graph Fourier
 128 coefficients of the original and reconstructed features on the Cora dataset, indicating that the higher
 129 the rate of missing features, the stronger the low-pass filtering effect.

130 The following results shows that the iterative scheme with $h = 1$ always converges and its steady
 131 state is equal to the closed form solution. Importantly, the solution does not depend on the initial
 132 values $\mathbf{x}_u^{(0)}$ given to the unknown features.

133 **Proposition 3.2.** Take any undirected and connected graph with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$,
 134 and normalised Adjacency $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$, with \mathbf{D} being the degree matrix of \mathbf{A} . Let $\mathbf{x} =$
 135 $\mathbf{x}^{(0)} \in \mathbf{R}^n$ be the initial feature vector and define the following recursive relation

$$\mathbf{x}^{(k)} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \tilde{\mathbf{A}}_{uk} & \tilde{\mathbf{A}}_{uu} \end{bmatrix} \mathbf{x}^{(k-1)}.$$

136 Then this recursion converges and the steady state is given to be

$$\lim_{n \rightarrow \infty} \mathbf{x}^{(n)} = \begin{bmatrix} \mathbf{x}_k \\ -\Delta_{kk}^{-1} \tilde{\mathbf{A}}_{uk} \mathbf{x}_k \end{bmatrix}.$$

137 Proof: See Appendix A.3.

138 **Feature Propagation Algorithm.** We can notice that the update in Equation 2 is equivalent to
 139 first multiplying the feature vector \mathbf{x} by the original diffusion matrix $\tilde{\mathbf{A}}$, and then resetting the
 140 known features to their true value. This gives us Algorithm 1, an extremely simple and scalable
 141 iterative algorithm to reconstruct the missing features on a graph, which we refer to as *Feature*
 142 *Propagation* (FP). While \mathbf{x}_u can be initialized to any value, in practice we initialize \mathbf{x}_u to zero
 143 and find 40 iterations to be enough to provide convergence for all datasets we experimented
 144 on. At each iteration, the diffusion occurs from the nodes with known features to the nodes with
 145 unknown features as well as among the nodes with unknown features.

Algorithm 1 Feature Propagation

```

1: Input: feature vector  $\mathbf{x}$ , diffusion matrix  $\tilde{\mathbf{A}}$ 
2:  $\mathbf{y} \leftarrow \mathbf{x}$ 
3: while  $\mathbf{x}$  has not converged do
4:    $\mathbf{x} \leftarrow \tilde{\mathbf{A}}\mathbf{x}$            ▷ Propagate features
5:    $\mathbf{x}_k \leftarrow \mathbf{y}_k$        ▷ Reset known features
6: end while

```

152 4 Related Work

153 **Label Propagation.** The proposed algorithm bears some similarity with Label Propagation [58]
 154 (LP), which predicts a class for each node by propagating the known labels in the graph. Differently
 155 from our setting of diffusion of continuous node features, they deal with discrete label classes directly,
 156 resulting in a different diffusion operator. However, the key difference between them lies in how they
 157 are used. Importantly, LP is used to directly perform node classification, taking into account only
 158 the graph structure and being unaware of node features. On the other hand, FP is used to reconstruct
 159 missing features, which are then fed into a downstream GNN classifier. FP allows a GNN model to
 160 effectively combine features and graph structures, even when most of the features are missing. Our
 161 experiments show that FP+GNN always outperforms LP, even in cases of extremely high rates of
 162 missing features, suggesting the effectiveness of FP. Also, the derived scheme is a special case of
 163 Neural Graph PDEs [8], which are in turn related to the iterative scheme presented in [56].

164 **Matrix completion.** Several optimization-based approaches [7, 25] as well as learning-based
 165 approaches [29, 32, 55] have been proposed to solve the matrix completion problem. However, they
 166 are unaware of the underlying graph structure. Graph matrix completion [27, 33, 37, 46] extends
 167 the above approaches to make use of an underlying graph. Similarly, Graph Signal Processing
 168 offers several methods to interpolate signals on graphs. [34] prove the necessary conditions for a
 169 graph signal to be recovered perfectly, and provide a corresponding algorithm. However, due to
 170 the optimisation problems involved, most above approaches are too computationally intensive and
 171 cannot scale to graphs with more than $\sim 1,000$ nodes. Moreover, the goal of all above approaches is
 172 to reconstruct the missing entries of the matrix, rather than solving a downstream task.

173 **Extending GNNs to missing node features.** SAT [10] consists of a Transformer-like model for
 174 feature reconstruction and a GNN model to solve the downstream task. GCNMF [44] adapts GCN [30]
 175 to the case of missing node features by representing the missing data with a Gaussian mixture model.
 176 PaGNN [26] is a GCN-like model which uses a partial message-passing scheme to only propagate
 177 observed features. While showing a reasonable performance for low rates of missing features, these
 178 methods suffer in regimes of high rates of missing features, and do not scale to large graphs.

179 **Other related GNN works.** Several papers investigate how to augment GNNs when no node
180 features are available [12], as well as investigating the performance of GNNs with random features [1,
181 39]. Dirichlet energy minimization has been widely used as a regularizer in several graph-related
182 tasks [49, 56, 59]. Discretization of continuous diffusion on graphs has already been explored in [8]
183 and [51]. Propagation on the graph has also been studied as a solution to the different problem of
184 node regression on multi-relational graphs [4]. Other methods have investigated propagating node
185 features [9, 18, 50], however not in the scenario of missing features. The boundary conditions given
186 by the available features in FP’s diffusion equation (enforced by resetting the known feature after
187 each iteration in the algorithm) is what makes it different from other propagation approaches and
188 makes it an effective solution to the missing features problem. While [9, 18, 50] assume to observe
189 all features, and then modify all features, FP assumes to observe only a subset of the features and
190 modifies only the unobserved ones.

191 5 Experiments and Discussion

192 **Datasets.** We evaluate on the task of node classification on several benchmark datasets: Cora,
193 Citeseer and PubMed [41], Amazon-Computers, Amazon-Photo [42] and OGBN-Arxiv [24]. To
194 test the scalability of our method, we also test it on OGBN-Products (2,449,029 nodes, 123,718,280
195 edges). We report dataset statistics in table 3 (Appendix).

196 **Baselines.** We compare to two strong feature-agnostic baselines: Label Propagation [58], which
197 only makes use of the graph structure by propagating labels on the graph, and Graph Positional
198 Encodings [15], which consist in computing the top k eigenvectors of the Laplacian matrix and
199 treating them as node features in input to a GNN. We additionally compare to feature-imputation
200 methods that are graph-agnostic, such as setting the missing features to 0 (Zero), a random value from
201 a standard Gaussian (Random), or the global mean of that feature over the graph (Global Mean)². We
202 also compare to a simple graph-based imputation baseline, which sets a missing feature to the mean
203 (of that same feature) over the neighbors of a node (Neighbor Mean). We additionally experiment with
204 MGCNN [33], a geometric graph completion method which learns how to reconstruct the missing
205 features by making use of the observed features and the graph structure. For all the above baselines,
206 as well as for our Feature Propagation, we experiment with both GCN [30] and GraphSage with
207 mean aggregator [23] as downstream GNNs. We also compare to recently state-of-the-art methods
208 for learning in the missing features setting (GCNMF [44] and PaGNN [26]). For GCNMF we use
209 the publicly available code.³ We could not find publicly available code for PaGNN so use our own
210 implementation for this comparison. We do not compare to other commonly used imputation based
211 methods such as VAE [29] or GAIN [55], nor to the Transformer-based method SAT [10], as they
212 have previously been shown to consistently underperform GCNMF and PaGNN [26, 44].

213 **Experimental Setup.** We report the mean and standard error of the test accuracy, computed over 10
214 runs, in all experiments. Each run has a different train/validation/test split (apart from OGBN datasets
215 where we use the provided splits) and mask of missing features⁴. The splits are generated at random
216 by assigning 20 nodes per class to the training set, 1500 nodes in total to the validation set and the
217 rest to the test set, similar to [31]. For a fair comparison, we use the same standard hyperparameters
218 for all methods across all experiments. We train using the Adam [28] optimizer with a learning rate
219 of 0.005 for a maximum of 10000 epochs, combined with early stopping with a patience of 200.
220 Downstream GNN models (as well as GCNMF and PaGNN) use 2 layers with a hidden dimension of
221 64 and a dropout rate of 0.5 for all datasets, apart from OGBN datasets where 3 layers and a hidden
222 dimension of 256 are used. For OGBN-Arxiv we also employ the Jumping Knowledge scheme [53]
223 with max aggregation. Feature Propagation uses 40 iterations to diffuse the features, as we found this
224 to be enough to reach convergence on all datasets. We want to emphasize that we did not perform
225 any hyperparameter tuning, and FP proved to perform consistently with any reasonable choice of
226 hyperparameters. We use neighbor sampling [23] when training on OGBN-Products. All experiments
227 are conducted on an AWS p3.16xlarge machine with 8 NVIDIA V100 GPUs with 16GB of memory
228 each, and took around 4 GPU days in total to perform.

²If a feature is not observed for any of the node’s neighbors, we set it to zero.

³<https://github.com/marblet/GCNmf>

⁴Each entry of the feature matrix is independently missing with a probability equal to the missing rate.

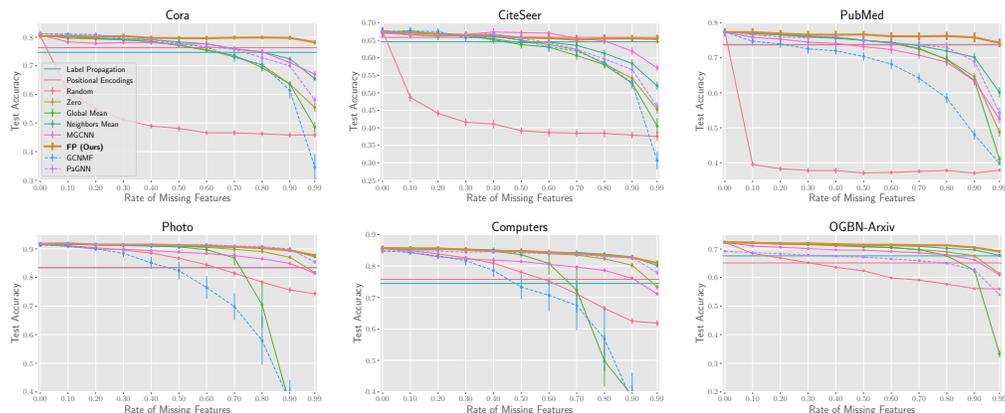


Figure 3: Test accuracy for varying rate of missing features on six common node-classification benchmarks. For methods that require a downstream GNNs, a 2-layer GCN [30] is used. On OGBN-Arxiv, GCNMF goes out-of-memory and is not reported.

Dataset	Full Features	50.0% Missing	90.0% Missing	99.0% Missing
Cora	80.39%	79.70%(-0.86%)	79.77%(-0.77%)	78.22%(-2.70%)
CiteSeer	67.48%	65.74%(-2.57%)	65.57%(-2.82%)	65.40%(-3.08%)
PubMed	77.36%	76.68%(-0.89%)	75.85%(-1.96%)	74.29%(-3.97%)
Photo	91.73%	91.29%(-0.48%)	89.48%(-2.46%)	87.73%(-4.36%)
Computers	85.65%	84.77%(-1.04%)	82.71%(-3.43%)	80.94%(-5.51%)
OGBN-Arxiv	72.22%	71.42%(-1.10%)	70.47%(-2.43%)	69.09%(-4.33%)
OGBN-Products	78.70%	77.16%(-1.96%)	75.94%(-3.51%)	74.94%(-4.78%)
Average	79.08%	78.11%(-1.27%)	77.11%(-2.48%)	75.80%(-4.10%)

Table 1: Performance of Feature Propagation (combined with a GCN model) for 50%, 90% and 99% of missing features, and relative drop compared to the performance of the same model when all features are present. On average, our method loses only 2.50% of relative accuracy with 90% of missing features, and 4.12% with 99% of missing features.

Dataset	GCNMF	PaGNN	Label Prop.	Pos. Enc.	FP (Ours)
Cora	34.54±2.07	58.03±0.57	74.68±0.36	76.33±0.26	78.22±0.32
CiteSeer	30.65±1.12	46.02±0.58	64.60±0.40	65.87±0.37	65.40±0.54
PubMed	39.80±0.25	54.25±0.70	73.81±0.56	73.70±0.29	74.29±0.55
Photo	29.64±2.78	85.41±0.28	83.45±0.94	83.45±0.26	87.73±0.27
Computers	30.74±1.95	77.91±0.33	74.48±0.61	75.77±0.47	80.94±0.37
OGBN-Arxiv	OOM	53.98±0.08	67.56±0.00	65.08±0.04	69.09±0.06
OGBN-Products	OOM	OOM	74.42±0.00	OOM	74.94±0.07

Table 2: Performance of GCNMF, PaGNN and FP(+GCN) with 99% of features missing, as well as Label Propagation and Positional Encodings (which are feature-agnostic). GCNMF and PaGNN perform respectively 58.33% and 21.25% worse in terms of relative accuracy in this scenario compared to when all the features are present. In comparison, FP has only a 4.12% drop.

229 **Node Classification Results.** Figure 3 shows the results for different rates of missing features
 230 (x-axis), when using GCN as a downstream GNN (results with GraphSAGE are reported in Figure 6
 231 of the Appendix). FP matches or outperforms other methods in all scenarios. Both GCNMF and
 232 PaGNN are consistently outperformed by the simple Neighbor Mean baseline. This is not completely
 233 unexpected, as Neighbor Mean can be seen as a first-order approximation of Feature Propagation,
 234 where only one step of propagation is performed (and with a slightly different normalization of the

diffusion operator). We elaborate on the relation between Neighbor Mean and Feature Propagation as well as on the results of the other baselines in Section A.5 of the Appendix. Interestingly, most methods perform extremely well up to 50% of missing features, suggesting that in general node features are redundant, as replacing half of them with zeros (*Zero* baseline) has little effect on the performance. The gap between methods opens up from around 60% of missing features, and is particularly large for extremely high rates of missing features (90% or 99%): FP is the only feature-aware method which is robust to these high rates on all datasets (see Table 2). Moreover, FP outperforms or matches Label Propagation and Positional Encodings on all datasets, even in the extreme case of 99% missing features. On some datasets, such as Cora, Photo, and Computers, the gap is especially significant. We conclude that reconstructing the missing features using FP is indeed useful for the downstream task. We highlight the surprising results that, on average, FP with 99% missing features performs only 4.12% worse (in relative accuracy terms) than the same GNN model used with no missing features, compared to 58.33% and 21.25% worse for GCNMF and PaGNN respectively.

Run-time analysis. Feature Propagation scales to extremely large graphs, as it only consists of repeated sparse-to-dense matrix multiplications. Moreover, it can be regarded as a pre-processing step, and performed only once, separately from training. In Figure 4 we compare the run-time to complete the training of the model for FP, PaGNN and GCNMF. The time for FP includes both the feature propagation step to reconstruct the missing features, as well as training of a downstream GCN model. FP is around 3x faster than PaGNN and GCNMF. The propagation step of FP takes only a fraction of the total running time, and the vast majority of the time is spent in training of the downstream model. The feature propagation step takes only $\sim 0.6s$ for Computers, $\sim 0.8s$ for OGBN-Arxiv and $\sim 10.5s$ for OGBN-Products using a single GPU. Both PaGNN and GCNMF go out-of-memory on OGBN-Products.

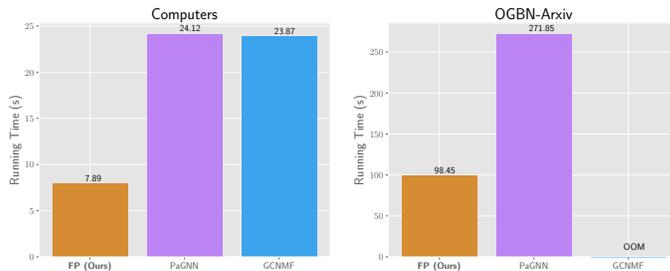


Figure 4: Run-time (in seconds) of FP, PaGNN and GCNMF. FP is 3x faster than both other methods. GCNMF goes out-of-memory (OOM) on OGBN-Arxiv.

6 Conclusion

We have introduced a novel approach for handling missing node features in graph-learning tasks. Our Feature Propagation model can be directly derived from energy minimization, and can be implemented as an efficient iterative algorithm where the features are multiplied by a diffusion matrix, before resetting the known features to their original value. Experiments on a number of datasets suggest that FP can reconstruct the missing features in a way that is useful for the downstream task, even when 99% of the features are missing. FP outperforms recently proposed methods by a significant margin on common benchmarks, while also being extremely scalable.

Limitations. While our method is designed for homophilic graphs, a more general learnable diffusion could be adopted to perform well in low homophily scenarios, as discussed in Section A.1. Feature Propagation is designed for graphs with only one node and edge type, however it could be extended to heterogenous graphs by having separate diffusions for different types of edges and nodes. Finally, Feature Propagation treats feature channels independently. To account for dependencies, diffusion with channel mixing should be used.

References

- [1] R. Abboud, İ. İ. Ceylan, M. Grohe, and T. Lukasiewicz. The surprising power of graph neural networks with random node initialization. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 2112–2118, 2021.

- 286 [2] S. Abu-El-Haija, B. Perozzi, A. Kapoor, H. Harutyunyan, N. Alipourfard, K. Lerman, G. V.
287 Steeg, and A. Galstyan. Mixhop: Higher-order graph convolution architectures via sparsified
288 neighborhood mixing. In *International Conference on Machine Learning (ICML)*, 2019.
- 289 [3] P. Battaglia, J. B. C. Hamrick, V. Bapst, A. Sanchez, V. Zambaldi, M. Malinowski, A. Tacchetti,
290 D. Raposo, A. Santoro, R. Faulkner, C. Gulcehre, F. Song, A. Ballard, J. Gilmer, G. E. Dahl,
291 A. Vaswani, K. Allen, C. Nash, V. J. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli,
292 M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu. Relational inductive biases, deep learning, and
293 graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- 294 [4] E. Bayram. Propagation on multi-relational graphs for node regression. In R. M. Benito,
295 C. Cherifi, H. Cherifi, E. Moro, L. M. Rocha, and M. Sales-Pardo, editors, *Complex Networks
296 & Their Applications X*, pages 155–167, Cham, 2022. Springer International Publishing.
- 297 [5] A. Berman and R. J. Plemmons. *Nonnegative Matrices in the Mathematical Sciences*. SIAM,
298 1994.
- 299 [6] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst. Geometric deep learning:
300 Going beyond Euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.
- 301 [7] E. J. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of
302 Computational mathematics*, 9(6):717–772, 2009.
- 303 [8] B. P. Chamberlain, J. R. Rowbottom, M. I. Gorinova, S. Webb, E. Rossi, and M. M. Bronstein.
304 GRAND: Graph neural diffusion. In *International Conference on Machine Learning (ICML)*,
305 2021.
- 306 [9] M. Chen, Z. Wei, B. Ding, Y. Li, Y. Yuan, X. D. Du, and J.-R. Wen. Scalable graph neural
307 networks via bidirectional propagation. 2020.
- 308 [10] X. Chen, S. Chen, J. Yao, H. Zheng, Y. Zhang, and I. Tsang. Learning on attribute-missing
309 graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP, 2020.
- 310 [11] F. R. K. Chung. *Spectral Graph Theory*. Number 92. American Mathematical Soc., 1997.
- 311 [12] H. Cui, Z. Lu, P. Li, and C. Yang. On positional and structural node features for graph neural
312 networks on non-attributed graphs. *International Workshop on Deep Learning on Graphs
313 (DLG-KDD)*, 2021.
- 314 [13] A. Darrow-Pinion, J. She, D. Wong, O. Lange, T. Hester, L. Perez, M. Nunkesser, S. Lee,
315 X. Guo, P. W. Battaglia, V. Gupta, A. Li, Z. Xu, A. Sanchez-Gonzalez, Y. Li, and P. Veličković.
316 Traffic Prediction with Graph Neural Networks in Google Maps. 2021.
- 317 [14] D. Duvenaud, D. Maclaurin, J. Aguilera-Iparraguirre, R. Gómez-Bombarelli, T. Hirzel,
318 A. Aspuru-Guzik, and R. P. Adams. Convolutional networks on graphs for learning molecu-
319 lar fingerprints. In *Proceedings of the 28th International Conference on Neural Information
320 Processing Systems*, pages 2224–2232, 2015.
- 321 [15] V. P. Dwivedi, C. K. Joshi, A. T. Luu, T. Laurent, Y. Bengio, and X. Bresson. Benchmarking
322 graph neural networks. *arXiv preprint arXiv:2003.00982*, 2020.
- 323 [16] M. Fey and J. E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR
324 Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- 325 [17] P. Gainza, F. Sverrisson, F. Monti, E. Rodolà, M. M. Bronstein, and B. E. Correia. Deciphering
326 interaction fingerprints from protein molecular surfaces. *Nature Methods*, 17(2):184–192, 2020.
- 327 [18] J. Gasteiger, A. Bojchevski, and S. Günnemann. Combining neural networks with personalized
328 pagerank for classification on graphs. In *International Conference on Learning Representations*,
329 2019.
- 330 [19] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing
331 for quantum chemistry. In *International conference on machine learning*, pages 1263–1272.
332 PMLR, 2017.

- 333 [20] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for
334 quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning*,
335 volume 70 of *Proceedings of Machine Learning Research*, pages 1263–1272, 2017.
- 336 [21] M. Gori, G. Monfardini, and F. Scarselli. A new model for learning in graph domains. In
337 *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 2,
338 pages 729–734. IEEE, 2005.
- 339 [22] A. Grover and J. Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings*
340 *of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*,
341 pages 855–864, 2016.
- 342 [23] W. L. Hamilton, R. Ying, and J. Leskovec. Inductive representation learning on large graphs. In
343 *Proceedings of the 31st International Conference on Neural Information Processing Systems*,
344 pages 1025–1035, 2017.
- 345 [24] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec. Open Graph
346 Benchmark: Datasets for machine learning on graphs. *arXiv preprint arXiv:2005.00687*, 2020.
- 347 [25] Y. Hu, Y. Koren, and C. Volinsky. Collaborative filtering for implicit feedback datasets. In *2008*
348 *Eighth IEEE International Conference on Data Mining*, pages 263–272, 2008.
- 349 [26] B. Jiang and Z. Zhang. Incomplete graph representation and learning via partial graph neural
350 networks. *arXiv preprint arXiv:2003.10130*, 2021.
- 351 [27] V. Kalofolias, X. Bresson, M. M. Bronstein, and P. Vandergheynst. Matrix completion on
352 graphs. *ArXiv preprint arXiv:1408.1717*, 2014.
- 353 [28] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *3rd International*
354 *Conference on Learning Representations, ICLR, 2015*.
- 355 [29] D. P. Kingma and M. Welling. Auto-encoding variational Bayes. In *International Conference*
356 *on Learning Representations, ICLR, 2014*.
- 357 [30] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks.
358 In *International Conference on Learning Representations, ICLR, 2017*.
- 359 [31] J. Klicpera, S. Weißenberger, and S. Günnemann. Diffusion improves graph learning. In
360 *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- 361 [32] X. Liu, X. Zhu, M. Li, L. Wang, E. Zhu, T. Liu, M. Kloft, D. Shen, J. Yin, and W. Gao. Multiple
362 kernel k -means with incomplete kernels. *IEEE Transactions on Pattern Analysis and Machine*
363 *Intelligence*, 42(5):1191–1204, 2020.
- 364 [33] F. Monti, M. M. Bronstein, and X. Bresson. Geometric matrix completion with recurrent
365 multi-graph neural networks. In *Proceedings of the 31st International Conference on Neural*
366 *Information Processing Systems, NIPS’17*, page 3700–3710, 2017.
- 367 [34] S. K. Narang, A. Gadde, and A. Ortega. Signal processing techniques for interpolation in
368 graph structured data. In *2013 IEEE International Conference on Acoustics, Speech and Signal*
369 *Processing*, pages 5445–5449, 2013.
- 370 [35] K. Oono and T. Suzuki. Graph neural networks exponentially lose expressive power for node
371 classification. In *International Conference on Learning Representations*, 2020.
- 372 [36] B. Perozzi, R. Al-Rfou, and S. Skiena. Deepwalk: Online learning of social representations. In
373 *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and*
374 *data mining*, pages 701–710, 2014.
- 375 [37] N. Rao, H.-F. Yu, P. K. Ravikumar, and I. S. Dhillon. Collaborative filtering with graph informa-
376 tion: Consistency and scalable methods. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and
377 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran
378 Associates, Inc., 2015.

- 379 [38] A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. Battaglia. Learning
380 to simulate complex physics with graph networks. In *International Conference on Machine*
381 *Learning (ICML)*, 2020.
- 382 [39] R. Sato, M. Yamada, and H. Kashima. Random features strengthen graph neural networks. In
383 *Proceedings of the 2021 SIAM International Conference on Data Mining, SDM*, pages 333–341.
384 SIAM, 2021.
- 385 [40] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural
386 network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
- 387 [41] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Gallagher, and T. Eliassi-Rad. Collective classifica-
388 tion in network data. *AI Magazine*, 29(3):93–106, 2008.
- 389 [42] O. Shchur, M. Mumme, A. Bojchevski, and S. Günnemann. Pitfalls of graph neural network
390 evaluation. *Relational Representation Learning Workshop, NeurIPS*, 2018.
- 391 [43] J. Shlomi, P. Battaglia, and J.-R. Vlimant. Graph neural networks in particle physics. *Machine*
392 *Learning: Science and Technology*, 2(2):021001, 2020.
- 393 [44] H. Taguchi, X. Liu, and T. Murata. Graph convolutional networks for graphs containing missing
394 features. *Future Generation Computer Systems*, 117:155–168, 2021.
- 395 [45] M. Thorpe, T. M. Nguyen, H. Xia, T. Strohmer, A. Bertozzi, S. Osher, and B. Wang. GRAND++:
396 Graph neural diffusion with a source term. In *International Conference on Learning Representations*,
397 2022.
- 398 [46] R. van den Berg, T. N. Kipf, and M. Welling. Graph convolutional matrix completion. *arXiv*
399 *preprint arXiv:1706.02263*, 2017.
- 400 [47] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph attention
401 networks. *International Conference on Learning Representations, ICLR*, 2018.
- 402 [48] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski,
403 P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman,
404 N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng,
405 E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero,
406 C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0
407 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature*
408 *Methods*, 17:261–272, 2020.
- 409 [49] J. Weston, F. Ratle, and R. Collobert. Deep learning via semi-supervised embedding. In
410 *Proceedings of the 25th International Conference on Machine Learning*, pages 1168–1175,
411 New York, NY, USA, 2008. Association for Computing Machinery.
- 412 [50] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger. Simplifying graph convolutional
413 networks. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International*
414 *Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*,
415 pages 6861–6871. PMLR, 09–15 Jun 2019.
- 416 [51] L.-P. Xhonneux, M. Qu, and J. Tang. Continuous graph neural networks. In *International*
417 *Conference on Machine Learning*, pages 10432–10441. PMLR, 2020.
- 418 [52] L.-P. A. C. Xhonneux, M. Qu, and J. Tang. Continuous graph neural networks. In *Proceedings*
419 *of the 37th International Conference on Machine Learning, ICML’20*. JMLR.org, 2020.
- 420 [53] K. Xu, C. Li, Y. Tian, T. Sonobe, K. ichi Kawarabayashi, and S. Jegelka. Representation learning
421 on graphs with jumping knowledge networks. In J. Dy and A. Krause, editors, *Proceedings of*
422 *the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine*
423 *Learning Research*, pages 5453–5462. PMLR, 10–15 Jul 2018.
- 424 [54] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec. Graph convolu-
425 tional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM*
426 *SIGKDD International Conference on Knowledge Discovery & Data Mining*, page 974–983.
427 Association for Computing Machinery, 2018.

- 428 [55] J. Yoon, J. Jordon, and M. van der Schaar. GAIN: Missing data imputation using generative
 429 adversarial nets. In *Proceedings of the 35th International Conference on Machine Learning*
 430 (*ICML*), volume 80 of *Proceedings of Machine Learning Research*, pages 5689–5698. PMLR,
 431 2018.
- 432 [56] D. Zhou and B. Schölkopf. A regularization framework for learning from graph data. In
 433 *Workshop on Statistical Relational Learning (ICML)*, 2004.
- 434 [57] J. Zhu, Y. Yan, L. Zhao, M. Heimann, L. Akoglu, and D. Koutra. Beyond homophily in graph
 435 neural networks: Current limitations and effective designs. *Advances in Neural Information*
 436 *Processing Systems (NeurIPS)*, 2020.
- 437 [58] X. Zhu and Z. Ghahramani. Learning from labeled and unlabeled data with label propagation.
 438 In *Technical Report CMU-CALD-02-107*, Carnegie Mellon University, 2002.
- 439 [59] X. Zhu, Z. Ghahramani, and J. Lafferty. Semi-supervised learning using gaussian fields and
 440 harmonic functions. In *Proceedings of the Twentieth International Conference on International*
 441 *Conference on Machine Learning*, ICML’03, page 912–919. AAAI Press, 2003.
- 442 [60] M. Zitnik, M. Agrawal, and J. Leskovec. Modeling polypharmacy side effects with graph
 443 convolutional networks. *Bioinformatics*, 34(13):i457–i466, 2018.

444 A Appendix

445 A.1 Algorithm Discussion

446 **Extension to Vector-Valued Features.** Algorithm 1 extends seamlessly to vector-valued features
 447 by simply replacing the feature vector \mathbf{x} with a $n \times d$ feature matrix \mathbf{X} , where d is the number of
 448 features. Multiplying the diffusion matrix \mathbf{A} by the feature matrix \mathbf{X} diffuses each feature channel
 449 independently. Interestingly, it would not be trivial to extend Equation 2 to vector-valued features,
 450 without noticing its equivalence with Algorithm 1, as each node could have different missing features,
 451 leading to different sub-matrices $\tilde{\mathbf{A}}_{uk}$ and $\tilde{\mathbf{A}}_{uu}$ for each feature channel.

452 **Learning.** One significant advantage of FP is that it can be easily combined with any graph
 453 learning model to generate predictions for the downstream task. Moreover, FP is not aimed at merely
 454 reconstructing the node features. Instead, by only reconstructing the lower frequency components
 455 of the signal, it is by design very well suited to be combined with GNNs, which are known to
 456 mainly leverage these lower frequency components [50]. Our approach is generic and can be used for
 457 any graph-related task for missing features, such as node classification, link prediction and graph
 458 classification. In this paper, we focus on node classification.

459 **Oversmoothing.** Figure 2 shows that the more features are missing, the smoother the reconstruction
 460 produced by FP is. Despite this, FP does not suffer from oversmoothing [35], a term used when node
 461 representations converge to similar values. Oversmoothing is caused by repeated diffusion and occurs
 462 widely when stacking more than a few layers of the most popular GNNs such as GCN [30], GAT [47]
 463 or SGC [50]. However, the boundary conditions in the Feature Propagation diffusion equation prevent
 464 the reconstructed features from becoming overly smooth, even when using an extremely high number
 465 of diffusion steps. This has also been studied by CGNN [52] and GRAND++ [45], which require
 466 soft boundary conditions in the form of a source term to prevent oversmoothing, although not in the
 467 context of missing features.

468 **When does Feature Propagation work?** Since FP can be interpreted as a low-pass filter that
 469 smoothes the features on the graph, we expect it to be suitable in the case of homophilic graph
 470 data (where neighbors tend to have similar attributes), and, conversely, to suffer in scenarios of low
 471 homophily. To verify this, we experiment on the synthetic dataset from [2], which consists of 10
 472 graphs with different levels of homophily. Figure 5 confirms our hypothesis: when the homophily
 473 is high, Feature Propagation with 99% of features missing performs similarly to the case when all
 474 the features are known. As the homophily decreases, the gap between the two widens to become
 475 extremely large in the case of zero homophily. In such scenarios, FP is only slightly better than

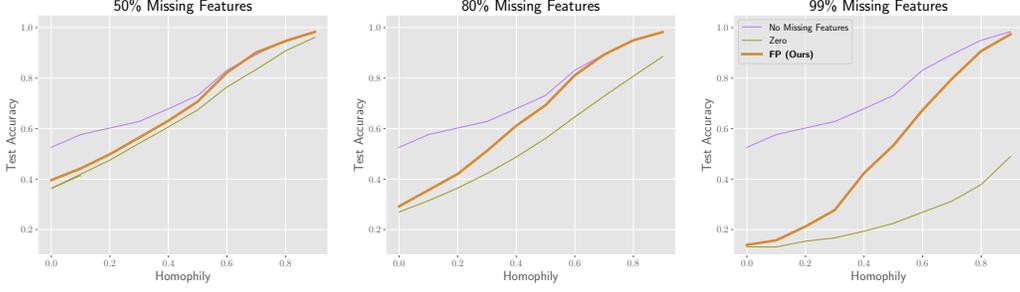


Figure 5: Test accuracy on the synthetic datasets from [2] with different levels of homophily. We use GraphSage as downstream model as it is preferable to GCN on low homophily data [57].

476 setting the missing features to zero (Zero baseline). This observation calls for a different kind of
 477 non-homogeneous diffusion dependent on the features that can potentially be made learnable for
 478 low-homophily data. We leave this as future work.

479 A.2 Closed-Form Solution for Harmonic Interpolation

480 Given the *Dirichlet energy* $\ell(\mathbf{x}, G) = \frac{1}{2} \mathbf{x}^\top \Delta \mathbf{x}$, we want to solve for missing features $\mathbf{x}_u =$
 481 $\text{argmin}_{\mathbf{x}_u} \ell$, leading to the optimality condition $\nabla_{\mathbf{x}_u} \ell = \mathbf{0}$. From Eq. 1 we find $\nabla_{\mathbf{x}_u} \ell = \mathbf{0}$ to
 482 be the solution of $\Delta_{uk} \mathbf{x}_k + \Delta_{uu} \mathbf{x}_u = \mathbf{0}$. The unique solution to this system of linear equations
 483 is $\mathbf{x}_u = -\Delta_{uu}^{-1} \Delta_{uk} \mathbf{x}_k$. We show this solution always exists by proving Δ_{uu} is non-singular
 484 (Proposition 3.1). The proof of this result follows from the following Lemma.

485 **Lemma A.1.** *Take any undirected and connected graph with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$, and*
 486 *normalised Adjacency $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$, with \mathbf{D} being the degree matrix of \mathbf{A} . Let $\tilde{\mathbf{A}}_{uu}$ be the*
 487 *bottom right submatrix of $\tilde{\mathbf{A}}$ where $1 \leq b < n$. Then $\rho(\tilde{\mathbf{A}}_{uu}) < 1$ where $\rho(\cdot)$ denotes spectral radius.*

488 *Proof.* Define

$$\tilde{\mathbf{A}}_{up} = \begin{bmatrix} \mathbf{0}_u & \mathbf{0}_{uk} \\ \mathbf{0}_{ku} & \tilde{\mathbf{A}}_{uu} \end{bmatrix},$$

489 to be the matrix equal to $\tilde{\mathbf{A}}_{uu}$ in the lower right $b \times b$ sub-matrix and padded with zero entries
 490 elsewhere. Clearly $\tilde{\mathbf{A}}_{up} \leq \tilde{\mathbf{A}}$ elementwise and $\tilde{\mathbf{A}}_{up} \neq \tilde{\mathbf{A}}$. Furthermore, $\tilde{\mathbf{A}}_{up} + \tilde{\mathbf{A}}$ represents
 491 an adjacency matrix of some strongly connected graph and is therefore irreducible [5, Theorem
 492 2.2.7]. These observations allow us to deduce that $\rho(\tilde{\mathbf{A}}_{up}) < \rho(\tilde{\mathbf{A}})$ [5, Corollary 2.1.5]. Note that
 493 $\rho(\tilde{\mathbf{A}}_{up}) = \rho(\tilde{\mathbf{A}}_{uu})$ as $\tilde{\mathbf{A}}_{up}$ and $\tilde{\mathbf{A}}_{uu}$ share the same non-zero eigenvalues. Furthermore, $\rho(\tilde{\mathbf{A}}) \leq 1$
 494 as we can write $\tilde{\mathbf{A}} = \mathbf{I} - \Delta$ and Δ is known to have eigenvalues in the range $[0, 2]$ [11]. Combining
 495 these inequalities gives the result $\rho(\tilde{\mathbf{A}}_{uu}) = \rho(\tilde{\mathbf{A}}_{up}) < \rho(\tilde{\mathbf{A}}) \leq 1$. \square

496 **Proposition A.2** (The sub-Laplacian matrix of a undirected connected graph is invertible). *Take*
 497 *any undirected, connected graph with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$, and its Laplacian $\Delta =$*
 498 *$\mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$, with \mathbf{D} being the degree matrix of \mathbf{A} . Then, for any principle sub-matrix*
 499 *$\mathbf{L}_u \in \mathbb{R}^{b \times b}$ of the Laplacian, where $1 \leq b < n$, \mathbf{L}_u is invertible.*

500 *Proof.* To prove Δ_{uu} is non-singular it is enough to show 0 is not an eigenvalue. Note that $\Delta_{uu} =$
 501 $\mathbf{I} - \tilde{\mathbf{A}}_{uu}$ so 0 is not an eigenvalue if and only if $\tilde{\mathbf{A}}_{uu}$ does not have an eigenvalue equal to 1, which
 502 follows from Lemma A.1. \square

503 A.3 Closed-Form Solution for the Euler scheme

504 **Proposition A.3.** *Take any undirected and connected graph with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$,*
 505 *and normalised Adjacency $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$, with \mathbf{D} being the degree matrix of \mathbf{A} . Let $\mathbf{x} =$*

Dataset	Nodes	Edges	Features	Classes
Cora	2,485	5,069	1,433	7
CiteSeer	2,120	3,679	3,703	6
PubMed	19,717	44,324	500	3
Photo	7,487	119,043	745	8
Computers	13,381	245,778	767	10
OGBN-Arxiv	169,343	1,166,243	128	40
OGBN-Products	2,449,029	123,718,280	100	47

Table 3: Dataset statistics.

506 $\mathbf{x}^{(0)} \in \mathbf{R}^n$ be the initial feature vector and define the following recursive relation

$$\mathbf{x}^{(k)} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \tilde{\mathbf{A}}_{uk} & \tilde{\mathbf{A}}_{uu} \end{bmatrix} \mathbf{x}^{(k-1)}.$$

507 Then this recursion converges and the steady state is given to be

$$\lim_{n \rightarrow \infty} \mathbf{x}^{(n)} = \begin{bmatrix} \mathbf{x}_k \\ -\Delta_{kk}^{-1} \tilde{\mathbf{A}}_{uk} \mathbf{x}_k \end{bmatrix}.$$

508 *Proof.* The recursive relation can be written in the following form

$$\begin{bmatrix} \mathbf{x}_k^{(k)} \\ \mathbf{x}_u^{(k)} \end{bmatrix} = \begin{bmatrix} \mathbf{I}_l & \mathbf{0}_{ku} \\ \tilde{\mathbf{A}}_{uk} & \tilde{\mathbf{A}}_{uu} \end{bmatrix} \begin{bmatrix} \mathbf{x}_k^{(k-1)} \\ \mathbf{x}_u^{(k-1)} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_k^{(k-1)} \\ \tilde{\mathbf{A}}_{uk} \mathbf{x}_k^{(k-1)} + \tilde{\mathbf{A}}_{uu} \mathbf{x}_u^{(k-1)} \end{bmatrix}.$$

509 The first l rows remain the same so we can write $\mathbf{x}_k^{(k)} = \mathbf{x}_k^{(k-1)} = \mathbf{x}_k$ and consider just the
510 convergence of the last u rows

$$\mathbf{x}_u^{(k-1)} = \tilde{\mathbf{A}}_{uk} \mathbf{x}_k + \tilde{\mathbf{A}}_{uu} \mathbf{x}_u^{(k-1)}.$$

511 We can look at the stationary behaviour by unrolling this recursion and taking the limit to find
512 stationary state

$$\lim_{n \rightarrow \infty} \mathbf{x}_u^{(n)} = \lim_{n \rightarrow \infty} \tilde{\mathbf{A}}_{uu}^n \mathbf{x}_u^{(0)} + \left(\sum_{i=1}^n \tilde{\mathbf{A}}_{uu}^{(i-1)} \right) \tilde{\mathbf{A}}_{uk} \mathbf{x}_k.$$

513 Using Lemma A.1 we find $\lim_{n \rightarrow \infty} \tilde{\mathbf{A}}_{uu}^n \mathbf{x}_u^{(0)} = \mathbf{0}$ and the geometric series converges giving us the
514 following limit

$$\lim_{n \rightarrow \infty} \mathbf{x}_u^{(n)} = \left(\mathbf{I}_u - \tilde{\mathbf{A}}_{uu} \right)^{-1} \tilde{\mathbf{A}}_{uk} \mathbf{x}_k = -\Delta_{kk}^{-1} \tilde{\mathbf{A}}_{uk} \mathbf{x}_k.$$

515

□

516 A.4 Baselines' Implementation and Tuning

517 **Label Propagation** We use the label propagation implementation provided in Pytorch-
518 Geometric [16]. Since the method is quite sensitive to the value of the α hyperpa-
519 rameter, we perform a gridsearch separately on each dataset over the following values:
520 [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.99].

521 **Positional Encodings** We compute the laplacian eigenvectors using SciPy [48] sparse eigenvectors
522 routines. We use the top twenty eigenvectors as positional encodings.

523 **MGCNN** We re-implement MGCNN [33] in Pytorch by taking inspiration from the authors' public
524 TensorFlow code⁵. For simplicity, we use the version of the model with only graph convolutional
525 layers and without an LSTM. For the matrix completion training process, we split the observed
526 features into 50% input data, 40% training targets and 10% validation data. Once the MGCNN
527 model is trained, we feed it the matrix with all the observed features to predict the whole feature
528 matrix. This reconstructed features matrix is then used as input for a downstream GNN (as for the
529 feature-imputation baselines).

⁵<https://github.com/fmonti/mgcnn>

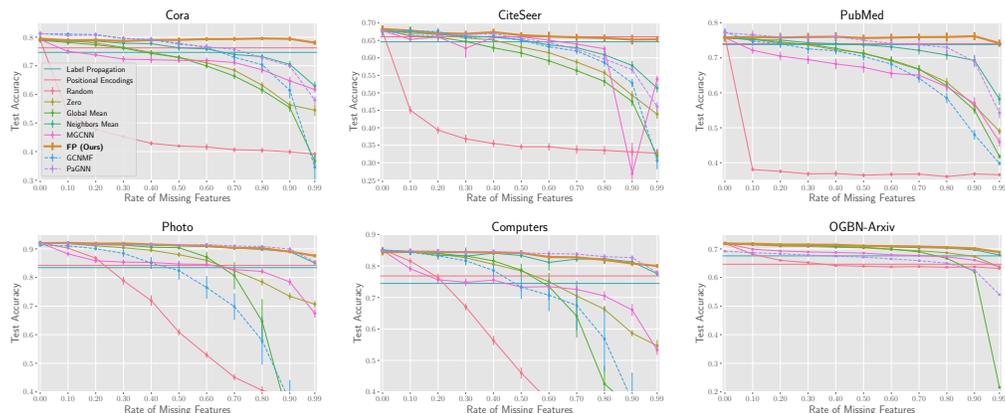


Figure 6: Test accuracy for varying rate of missing features on six common node-classification benchmarks. For methods that require a downstream GNNs, a 2-layer GraphSAGE [23] is used. On OGBN-Arxiv, GCNMF goes out-of-memory and is not reported.

530 A.5 Discussion Over Baselines' Performance

531 **Neighborhood Averaging** As for some intuition to why the simple Neighborhood Averaging per-
 532 forms competitively, let us assume to have a single feature channel and this feature to be homophilous
 533 over the graph. When a node has enough neighbors, the average of their features is a good estimate
 534 for the feature of the given node. However, as the rate of missing features increases, the feature may
 535 be present for only a few neighbors (or none at all), causing the estimate to have a higher variance.
 536 On the other hand, Feature Propagation allows information to travel longer distances in the graph by
 537 repeatedly multiplying by the diffusion matrix. Even if we do not observe the feature for any of a
 538 node's neighbors, it is still possible to estimate it from nodes further away in the graph. This can be
 539 observed empirically: the gap between Neighborhood Averaging and Feature Propagation becomes
 540 increasingly significant for higher rates of missing features.

541 **Zero vs Random** In models such as GCN and GraphSage, where node embeddings are computed
 542 as (weighted) average of neighbors embeddings, the effect of the Zero baseline is simply to reduce
 543 the norm of the average embeddings of all nodes (since all nodes have the same expected proportion
 544 of neighbors with missing features). On the other hand, the Random baseline corrupts this weighted
 545 average. More generally, while for a GNN model it could be relatively easy to learn to ignore features
 546 set to zero, and only focus on known (non-zero) features, it would be basically impossible for the
 547 model to do the same when setting the missing features to a random value.

548 However, we find Random to perform better than Zero when all features are missing. This is in
 549 line with findings in the literature [1, 39], where Random features have been shown to work well
 550 in conjunction with GNNs as they act as signatures for the nodes. On the other hand, if all nodes
 551 have all zero vectors, it becomes basically impossible to distinguish them. After applying a GNN, all
 552 nodes will still have very similar embeddings and the task performance will be close to a random
 553 guess.