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

Anonymous Author(s)
Anonymous Affiliation
Anonymous Email

Abstract

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

16 1 Introduction

17 Graph Neural Networks (GNNs) [1–6] have been successful on a broad range of problems and in a
18 variety of fields [7–13]. GNNs typically operate by a message-passing mechanism [14, 15], where
19 at each layer, nodes send their feature representations (“messages”) to their neighbors. The feature
20 representation of each node is initialized to their original features, and is updated by repeatedly aggregating
21 incoming messages from neighbors. Being able to combine the topological information with
22 feature information is what distinguishes GNNs from other purely topological learning approaches
23 such as random walks [16, 17] or label propagation [18], and 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 [19–21] can be used to fill the missing values of the feature matrix,
33 they are unaware of the underlying graph structure. Graph Signal Processing, a field attempting to
34 generalize classical Fourier analysis to graphs, offers several methods that reconstruct signals on
35 graphs [22]. However, they do not scale beyond graphs with a few thousand nodes, making them
36 infeasible for practical applications. More recently, SAT [23], GCNMF [24] and PaGNN [25] 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.

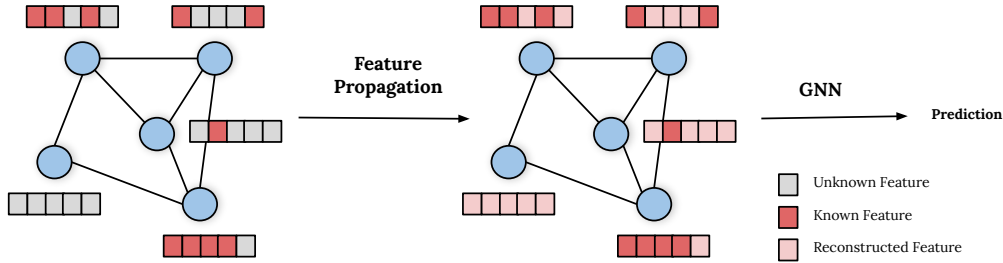


Figure 1: A diagram illustrating our Feature Propagation framework. On the left, a graph with missing node features. In the initial reconstruction step, Feature Propagation reconstructs the missing features by iteratively diffusing the known features in the graph. Subsequently, the graph and the reconstructed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

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
 43 machine learning tasks. The framework consists of an initial diffusion-based feature reconstruction
 44 step followed by a downstream GNN. The reconstruction step is based on Dirichlet energy
 45 minimization, which leads to a diffusion-type differential equation on the graph. Discretization of
 46 this differential equation leads to a very simple, fast, and scalable iterative algorithm which we call
 47 Feature 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.
- 54 • **Robust to high rates of missing features:** FP can withstand surprisingly high rates of missing
 55 features. In our experiment, we observe on average around 4% relative accuracy drop when up to
 56 99% of the features are missing. In comparison, GCNMF and PaGNN have an average drop of
 57 53.33% and 21.25% respectively. This finding has important implications especially in scenarios
 58 where the cost of sampling (observing features on nodes) is high or sampling is not possible
 59 altogether.
- 60 • **Generic:** FP can be combined with any GNN model to solve the downstream task; in contrast,
 61 GCNMF and PaGNN are specific GCN-type models.
- 62 • **Fast and Scalable:** FP takes only around 10 seconds for the reconstruction step on OGBN-
 63 Products (a graph with ~ 2.5 M nodes and ~ 123 M edges) on a single GPU. GCNMF and PaGNN
 64 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 $\Delta = \mathbf{I} - \tilde{\mathbf{A}}$, 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.

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

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

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} .

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.1, 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
 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 $\Delta =$*

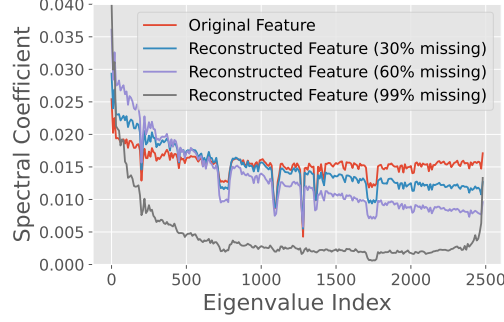


Figure 2: Graph Fourier transform magnitudes of the original Cora features (red) and those reconstructed by FP for varying rates of missing rates (we take the average over feature channels). Since FP minimizes the Dirichlet energy, it can be interpreted as a low-pass filter, which is stronger for a higher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

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.1. 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.2.

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

158 **Learning.** One significant advantage of FP is that it can be easily combined with any graph
 159 learning model to generate predictions for the downstream task. Moreover, FP is not aimed at merely
 160 reconstructing the node features. Instead, by only reconstructing the lower frequency components
 161 of the signal, it is by design very well suited to be combined with GNNs, which are known to
 162 mainly leverage these lower frequency components [26]. Our approach is generic and can be used for
 163 any graph-related task for missing features, such as node classification, link prediction and graph
 164 classification. In this paper, we focus on node classification.

165 **Oversmoothing.** Figure 2 shows that the more features are missing, the smoother the reconstruction
 166 produced by FP is. Despite this, FP does not suffer from oversmoothing [27], a term used when node
 167 representations converge to similar values. Oversmoothing is caused by repeated diffusion and occurs
 168 widely when stacking more than a few layers of the most popular GNNs such as GCN [3], GAT [5] or
 169 SGC [26]. However, the boundary conditions in the Feature Propagation diffusion equation prevent
 170 the reconstructed features from becoming overly smooth, even when using an extremely high number
 171 of diffusion steps. This has also been studied by CGNN [28] and GRAND++ [29], which require
 172 soft boundary conditions in the form of a source term to prevent oversmoothing, although not in the
 173 context of missing features.

174 **4 Related Work**

175 **Label Propagation.** The proposed algorithm bears some similarity with Label Propagation [18]
 176 (LP), which predicts a class for each node by propagating the known labels in the graph. Differently
 177 from our setting of diffusion of continuous node features, they deal with discrete label classes directly,
 178 resulting in a different diffusion operator. However, the key difference between them lies in how they
 179 are used. Importantly, LP is used to directly perform node classification, taking into account only
 180 the graph structure and being unaware of node features. On the other hand, FP is used to reconstruct
 181 missing features, which are then fed into a downstream GNN classifier. FP allows a GNN model to
 182 effectively combine features and graph structures, even when most of the features are missing. Our
 183 experiments show that FP+GNN always outperforms LP, even in cases of extremely high rates of
 184 missing features, suggesting the effectiveness of FP. Also, the derived scheme is a special case of
 185 Neural Graph PDEs [30], which are in turn related to the iterative scheme presented in [31].

186 **Matrix completion.** Several optimization-based approaches [32, 33] as well as learning-based
 187 approaches [19–21] have been proposed to solve the matrix completion problem. However, they
 188 are unaware of the underlying graph structure. Graph matrix completion [34–37] extends the above
 189 approaches to make use of an underlying graph. Similarly, Graph Signal Processing offers several
 190 methods to interpolate signals on graphs. [22] prove the necessary conditions for a graph signal to
 191 be recovered perfectly, and provide a corresponding algorithm. However, due to the optimisation
 192 problems involved, most above approaches are too computationally intensive and cannot scale to
 193 graphs with more than $\sim 1,000$ nodes. Moreover, the goal of all above approaches is to reconstruct
 194 the missing entries of the matrix, rather than solving a downstream task.

195 **Extending GNNs to missing node features.** SAT [23] consists of a Transformer-like model for
 196 feature reconstruction and a GNN model to solve the downstream task. GCNMF [24] adapts GCN [3]
 197 to the case of missing node features by representing the missing data with a Gaussian mixture model.
 198 PaGNN [25] is a GCN-like model which uses a partial message-passing scheme to only propagate
 199 observed features. While showing a reasonable performance for low rates of missing features, these
 200 methods suffer in regimes of high rates of missing features, and do not scale to large graphs.

201 **Other related GNN works.** Several papers investigate how to augment GNNs when no node
 202 features are available [38], as well as investigating the performance of GNNs with random features [39,
 203 40]. Dirichlet energy minimization has been widely used as a regularizer in several graph-related
 204 tasks [31, 41, 42]. Discretization of continuous diffusion on graphs has already been explored in [30]
 205 and [43]. Propagation on the graph has also been studied as a solution to the different problem of
 206 node regression on multi-relational graphs [44]. Other methods have investigated propagating node
 207 features [26, 45, 46], however not in the scenario of missing features. The boundary conditions given
 208 by the available features in FP’s diffusion equation (enforced by resetting the known feature after
 209 each iteration in the algorithm) is what makes it different from other propagation approaches and
 210 makes it an effective solution to the missing features problem. While [26, 45, 46] assume to observe
 211 all features, and then modify all features, FP assumes to observe only a subset of the features and
 212 modifies only the unobserved ones.

213 5 Experiments and Discussion

214 **Datasets.** We evaluate on the task of node classification on several benchmark datasets: Cora,
 215 Citeseer and PubMed [47], Amazon-Computers, Amazon-Photo [48] and OGBN-Arxiv [49]. To
 216 test the scalability of our method, we also test it on OGBN-Products (2,449,029 nodes, 123,718,280
 217 edges). We report dataset statistics in table 3 (Appendix).

218 **Baselines.** We compare to two strong feature-agnostic baselines: Label Propagation [18], which
 219 only makes use of the graph structure by propagating labels on the graph, and Graph Positional
 220 Encodings [50], which consist in computing the top k eigenvectors of the Laplacian matrix and
 221 treating them as node features in input to a GNN. We additionally compare to feature-imputation
 222 methods that are graph-agnostic, such as setting the missing features to 0 (Zero), a random value from
 223 a standard Gaussian (Random), or the global mean of that feature over the graph (Global Mean)². We
 224 also compare to a simple graph-based imputation baseline, which sets a missing feature to the mean
 225 (of that same feature) over the neighbors of a node (Neighbor Mean). We additionally experiment with
 226 MGCNN [36], a geometric graph completion method which learns how to reconstruct the missing
 227 features by making use of the observed features and the graph structure. For all the above baselines,
 228 as well as for our Feature Propagation, we experiment with both GCN [3] and GraphSage with mean
 229 aggregator [51] as downstream GNNs. We also compare to recently state-of-the-art methods for
 230 learning in the missing features setting (GCNMF [24] and PaGNN [25]). For GCNMF we use the
 231 publicly available code.³ We could not find publicly available code for PaGNN so use our own
 232 implementation for this comparison. We do not compare to other commonly used imputation based
 233 methods such as VAE [21] or GAIN [20], nor to the Transformer-based method SAT [23], as they
 234 have previously been shown to consistently underperform GCNMF and PaGNN [24, 25].

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

³<https://github.com/marblet/GCNmf>

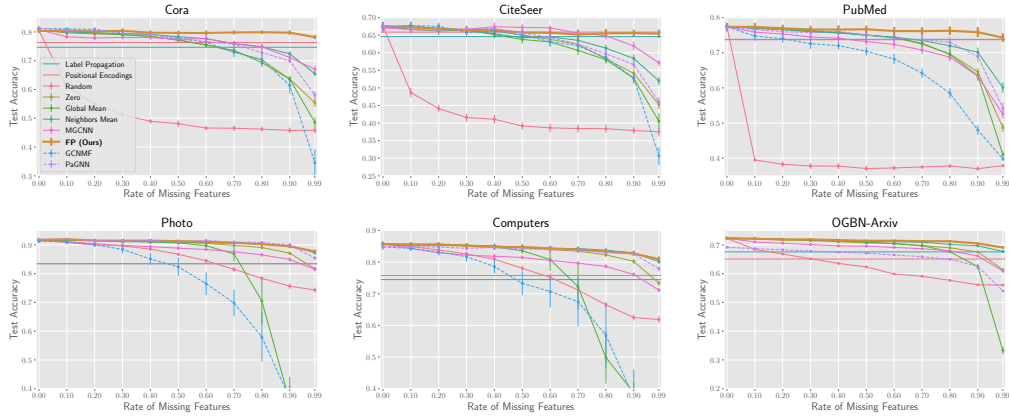


Figure 3: Test accuracy for varying rate of missing features on six common node-classification benchmarks. For methods that require a downstream GNNs, a 2-layer GCN [3] is used. On OGBN-Arxiv, GCNMF goes out-of-memory and is not reported.

235 **Experimental Setup.** We report the mean and standard error of the test accuracy, computed over 10
 236 runs, in all experiments. Each run has a different train/validation/test split (apart from OGBN datasets
 237 where we use the provided splits) and mask of missing features⁴. The splits are generated at random
 238 by assigning 20 nodes per class to the training set, 1500 nodes in total to the validation set and the
 239 rest to the test set, similar to [52]. For a fair comparison, we use the same standard hyperparameters
 240 for all methods across all experiments. We train using the Adam [53] optimizer with a learning rate
 241 of 0.005 for a maximum of 10000 epochs, combined with early stopping with a patience of 200.
 242 Downstream GNN models (as well as GCNMF and PaGNN) use 2 layers with a hidden dimension of
 243 64 and a dropout rate of 0.5 for all datasets, apart from OGBN datasets where 3 layers and a hidden
 244 dimension of 256 are used. For OGBN-Arxiv we also employ the Jumping Knowledge scheme [54]
 245 with max aggregation. Feature Propagation uses 40 iterations to diffuse the features, as we found this
 246 to be enough to reach convergence on all datasets. We want to emphasize that we did not perform
 247 any hyperparameter tuning, and FP proved to perform consistently with any reasonable choice of
 248 hyperparameters. We use neighbor sampling [51] when training on OGBN-Products. All experiments
 249 are conducted on an AWS p3.16xlarge machine with 8 NVIDIA V100 GPUs with 16GB of memory
 250 each, and took around 4 GPU days in total to perform.

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.

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

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.

251 **Node Classification Results.** Figure 3 shows the results for different rates of missing features
 252 (x-axis), when using GCN as a downstream GNN (results with GraphSAGE are reported in Figure 6
 253 of the Appendix). FP matches or outperforms other methods in all scenarios. Both GCNMF and
 254 PaGNN are consistently outperformed by the simple Neighbor Mean baseline. This is not completely
 255 unexpected, as Neighbor Mean can be seen as a first-order approximation of Feature Propagation,
 256 where only one step of propagation is performed (and with a slightly different normalization of the
 257 diffusion operator). We elaborate on the relation between Neighbor Mean and Feature Propagation
 258 as well as on the results of the other baselines in Section A.4 of the Appendix. Interestingly, most
 259 methods perform extremely well up to 50% of missing features, suggesting that in general node
 260 features are redundant, as replacing half of them with zeros (*Zero* baseline) has little effect on
 261 the performance. The gap between methods opens up from around 60% of missing features, and
 262 is particularly large for extremely high rates of missing features (90% or 99%): FP is the only
 263 feature-aware method which is robust to these high rates on all datasets (see Table 2). Moreover,
 264 FP outperforms or matches Label Propagation and Positional Encodings on all datasets, even in the
 265 extreme case of 99% missing features. On some datasets, such as Cora, Photo, and Computers, the
 266 gap is especially significant. We conclude that reconstructing the missing features using FP is indeed
 267 useful for the downstream task. We highlight the surprising results that, on average, FP with 99%
 268 missing features performs only 4.12% worse (in relative accuracy terms) than the same GNN model
 269 used with no missing features, compared to 58.33% and 21.25% worse for GCNMF and PaGNN
 270 respectively.

271 **Run-time analysis.** Feature
 272 Propagation scales to extremely
 273 large graphs, as it only consists
 274 of repeated sparse-to-dense matrix
 275 multiplications. Moreover,
 276 it can be regarded as a pre-
 277 processing step, and performed
 278 only once, separately from
 279 training. In Figure 4 we compare
 280 the run-time to complete the
 281 training of the model for FP,
 282 PaGNN and GCNMF. The time
 283 for FP includes both the feature
 284 propagation step to reconstruct
 285 the missing features, as well as
 286 training of a downstream GCN model. FP is around 3x faster than PaGNN and GCNMF. The
 287 propagation step of FP takes only a fraction of the total running time, and the vast majority of the
 288 time is spent in training of the downstream model. The feature propagation step takes only ~0.6s
 289 for Computers, ~0.8s for OGBN-Arxiv and ~10.5s for OGBN-Products using a single GPU. Both
 290 PaGNN and GCNMF go out-of-memory on OGBN-Products.

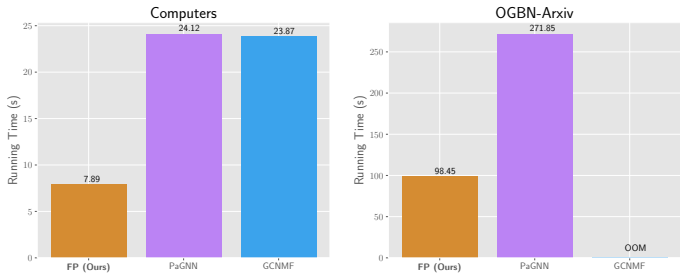


Figure 4: Run-time (in seconds) of FP, PaGNN and GCNMF. FP is 3x faster than both other methods. GCNMF goes out-of-memory (OOM) on OGBN-Arxiv.

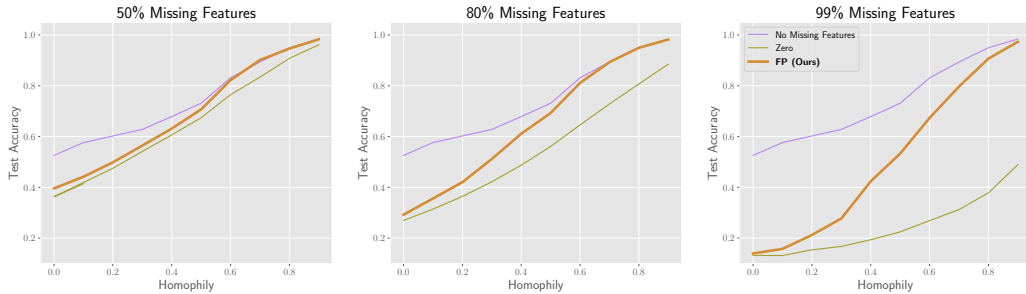


Figure 5: Test accuracy on the synthetic datasets from [55] with different levels of homophily. We use GraphSage as downstream model as it is preferable to GCN on low homophily data [56].

291 **When does Feature Propagation work?** Since FP can be interpreted as a low-pass filter that
 292 smoothes the features on the graph, we expect it to be suitable in the case of homophilic graph
 293 data (where neighbors tend to have similar attributes), and, conversely, to suffer in scenarios of low
 294 homophily. To verify this, we experiment on the synthetic dataset from [55], which consists of 10
 295 graphs with different levels of homophily. Figure 5 confirms our hypothesis: when the homophily
 296 is high, Feature Propagation with 99% of features missing performs similarly to the case when all
 297 the features are known. As the homophily decreases, the gap between the two widens to become
 298 extremely large in the case of zero homophily. In such scenarios, FP is only slightly better than
 299 setting the missing features to zero (Zero baseline). This observation calls for a different kind of
 300 non-homogeneous diffusion dependent on the features that can potentially be made learnable for
 301 low-homophily data. We leave this as future work.

302 6 Conclusion

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

310 **Limitations.** While our method is designed for homophilic graphs, a more general learnable
 311 diffusion could be adopted to perform well in low homophily scenarios, as discussed in Section 5.
 312 Feature Propagation is designed for graphs with only one node and edge type, however it could be
 313 extended to heterogeneous graphs by having separate diffusions for different types of edges and nodes.
 314 Finally, Feature Propagation treats feature channels independently. To account for dependencies,
 315 diffusion with channel mixing should be used.

316 **Societal Impact.** Our work is aimed at improving the performance of Graph Neural Networks.
 317 While we believe that nothing in our work raises specific ethical concerns, the recent broad adoption
 318 of GNNs in industrial applications opens the possibility to the misuse of such methods with potentially
 319 detrimental societal impact.

320 References

- 321 [1] Marco Gori, Gabriele Monfardini, and Franco Scarselli. A new model for learning in graph
 322 domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*,
 323 volume 2, pages 729–734. IEEE, 2005.
- 324 [2] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini.
 325 The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.

- 326 [3] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
327 networks. In *International Conference on Learning Representations, ICLR*, 2017.
- 328 [4] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
329 message passing for quantum chemistry. In *International conference on machine learning*,
330 pages 1263–1272. PMLR, 2017.
- 331 [5] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
332 Bengio. Graph attention networks. *International Conference on Learning Representations*,
333 *ICLR*, 2018.
- 334 [6] Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst.
335 Geometric deep learning: Going beyond Euclidean data. *IEEE Signal Processing Magazine*, 34
336 (4):18–42, 2017.
- 337 [7] David Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Gómez-Bombarelli,
338 Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P. Adams. Convolutional networks on graphs
339 for learning molecular fingerprints. In *Proceedings of the 28th International Conference on*
340 *Neural Information Processing Systems*, pages 2224–2232, 2015.
- 341 [8] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure
342 Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *Pro-*
343 *ceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data*
344 *Mining*, page 974–983. Association for Computing Machinery, 2018. ISBN 9781450355520.
- 345 [9] Marinka Zitnik, Monica Agrawal, and Jure Leskovec. Modeling polypharmacy side effects with
346 graph convolutional networks. *Bioinformatics*, 34(13):i457–i466, 2018.
- 347 [10] Pablo Gainza, Freyr Sverrisson, Federico Monti, Emanuele Rodolà, Michael M. Bronstein, and
348 Bruno E. Correia. Deciphering interaction fingerprints from protein molecular surfaces. *Nature*
349 *Methods*, 17(2):184–192, 2020.
- 350 [11] Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec, and
351 Peter Battaglia. Learning to simulate complex physics with graph networks. In *International*
352 *Conference on Machine Learning (ICML)*, 2020.
- 353 [12] Jonathan Shlomi, Peter Battaglia, and Jean-Roch Vlimant. Graph neural networks in particle
354 physics. *Machine Learning: Science and Technology*, 2(2):021001, 2020.
- 355 [13] Austin Derrow-Pinion, Jennifer She, David Wong, Oliver Lange, Todd Hester, Luis Perez, Marc
356 Nunkesser, Seongjae Lee, Xueying Guo, Peter W Battaglia, Vishal Gupta, Ang Li, Zhongwen
357 Xu, Alvaro Sanchez-Gonzalez, Yujia Li, and Petar Veličković. Traffic Prediction with Graph
358 Neural Networks in Google Maps. 2021.
- 359 [14] Peter Battaglia, Jessica Blake Chandler Hamrick, Victor Bapst, Alvaro Sanchez, Vinicius Zam-
360 baldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner,
361 Caglar Gulcehre, Francis Song, Andy Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani,
362 Kelsey Allen, Charles Nash, Victoria Jayne Langston, Chris Dyer, Nicolas Heess, Daan Wierstra,
363 Pushmeet Kohli, Matt Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational
364 inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- 365 [15] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl.
366 Neural message passing for quantum chemistry. In *Proceedings of the 34th International*
367 *Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*,
368 pages 1263–1272, 2017.
- 369 [16] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social repre-
370 sentations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge*
371 *discovery and data mining*, pages 701–710, 2014.
- 372 [17] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In
373 *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and*
374 *data mining*, pages 855–864, 2016.
- 375 [18] Xiaojin Zhu and Zoubin Ghahramani. Learning from labeled and unlabeled data with label
376 propagation. In *Technical Report CMU-CALD-02-107*, Carnegie Mellon University, 2002.
- 377 [19] Xinwang Liu, Xinzhong Zhu, Miaomiao Li, Lei Wang, En Zhu, Tongliang Liu, Marius Kloft,
378 Dinggang Shen, Jianping Yin, and Wen Gao. Multiple kernel k -means with incomplete kernels.
379 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(5):1191–1204, 2020.

- 380 [20] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. GAIN: Missing data imputation using
381 generative adversarial nets. In *Proceedings of the 35th International Conference on Machine*
382 *Learning (ICML)*, volume 80 of *Proceedings of Machine Learning Research*, pages 5689–5698.
383 PMLR, 2018.
- 384 [21] Diederik P. Kingma and Max Welling. Auto-encoding variational Bayes. In *International*
385 *Conference on Learning Representations, ICLR*, 2014.
- 386 [22] Sunil K. Narang, Akshay Gadde, and Antonio Ortega. Signal processing techniques for
387 interpolation in graph structured data. In *2013 IEEE International Conference on Acoustics,*
388 *Speech and Signal Processing*, pages 5445–5449, 2013.
- 389 [23] Xu Chen, Siheng Chen, Jiangchao Yao, Huangjie Zheng, Ya Zhang, and Ivor Tsang. Learning
390 on attribute-missing graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
391 PP, 2020.
- 392 [24] Hibiki Taguchi, Xin Liu, and Tsuyoshi Murata. Graph convolutional networks for graphs
393 containing missing features. *Future Generation Computer Systems*, 117:155–168, 2021.
- 394 [25] Bo Jiang and Ziyang Zhang. Incomplete graph representation and learning via partial graph
395 neural networks. *arXiv preprint arXiv:2003.10130*, 2021.
- 396 [26] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger.
397 Simplifying graph convolutional networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov,
398 editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of
399 *Proceedings of Machine Learning Research*, pages 6861–6871. PMLR, 09–15 Jun 2019.
- 400 [27] Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for
401 node classification. In *International Conference on Learning Representations*, 2020. URL
402 <https://openreview.net/forum?id=S1ld02EFPr>.
- 403 [28] Louis-Pascal A. C. Xhonneux, Meng Qu, and Jian Tang. Continuous graph neural networks. In
404 *Proceedings of the 37th International Conference on Machine Learning, ICML'20*. JMLR.org,
405 2020.
- 406 [29] Matthew Thorpe, Tan Minh Nguyen, Hedi Xia, Thomas Strohmer, Andrea Bertozzi, Stanley
407 Osher, and Bao Wang. GRAND++: Graph neural diffusion with a source term. In *International*
408 *Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=EMxu-dzvJk)
409 [id=EMxu-dzvJk](https://openreview.net/forum?id=EMxu-dzvJk).
- 410 [30] Benjamin Paul Chamberlain, James R. Rowbottom, Maria I. Gorinova, Stefan Webb, Emanuele
411 Rossi, and Michael M. Bronstein. GRAND: Graph neural diffusion. In *International Conference*
412 *on Machine Learning (ICML)*, 2021.
- 413 [31] Dengyong Zhou and Bernhard Schölkopf. A regularization framework for learning from graph
414 data. In *Workshop on Statistical Relational Learning (ICML)*, 2004.
- 415 [32] Emmanuel J. Candès and Benjamin Recht. Exact matrix completion via convex optimization.
416 *Foundations of Computational mathematics*, 9(6):717–772, 2009.
- 417 [33] Yifan Hu, Yehuda Koren, and Chris Volinsky. Collaborative filtering for implicit feedback
418 datasets. In *2008 Eighth IEEE International Conference on Data Mining*, pages 263–272, 2008.
- 419 [34] Vassilis Kalofolias, Xavier Bresson, Michael M. Bronstein, and Pierre Vandergheynst. Matrix
420 completion on graphs. *ArXiv preprint arXiv:1408.1717*, 2014.
- 421 [35] Rianne van den Berg, Thomas N Kipf, and Max Welling. Graph convolutional matrix completion.
422 *arXiv preprint arXiv:1706.02263*, 2017.
- 423 [36] Federico Monti, Michael M. Bronstein, and Xavier Bresson. Geometric matrix completion
424 with recurrent multi-graph neural networks. In *Proceedings of the 31st International Con-*
425 *ference on Neural Information Processing Systems, NIPS'17*, page 3700–3710, 2017. ISBN
426 9781510860964.
- 427 [37] Nikhil Rao, Hsiang-Fu Yu, Pradeep K Ravikumar, and Inderjit S Dhillon. Collaborative
428 filtering with graph information: Consistency and scalable methods. In C. Cortes, N. Lawrence,
429 D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing*
430 *Systems*, volume 28. Curran Associates, Inc., 2015. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper/2015/file/f4573fc71c731d5c362f0d7860945b88-Paper.pdf)
431 [cc/paper/2015/file/f4573fc71c731d5c362f0d7860945b88-Paper.pdf](https://proceedings.neurips.cc/paper/2015/file/f4573fc71c731d5c362f0d7860945b88-Paper.pdf).

- 432 [38] Hejie Cui, Zijie Lu, Pan Li, and Carl Yang. On positional and structural node features for graph
433 neural networks on non-attributed graphs. *International Workshop on Deep Learning on Graphs*
434 (*DLG-KDD*), 2021.
- 435 [39] Ryoma Sato, Makoto Yamada, and Hisashi Kashima. Random features strengthen graph neural
436 networks. In *Proceedings of the 2021 SIAM International Conference on Data Mining, SDM*,
437 pages 333–341. SIAM, 2021.
- 438 [40] Ralph Abboud, İsmail İlkan Ceylan, Martin Grohe, and Thomas Lukasiewicz. The surprising
439 power of graph neural networks with random node initialization. In *Proceedings of the Thirtieth*
440 *International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 2112–2118, 2021.
- 441 [41] Xiaojin Zhu, Zoubin Ghahramani, and John Lafferty. Semi-supervised learning using gaussian
442 fields and harmonic functions. In *Proceedings of the Twentieth International Conference on*
443 *International Conference on Machine Learning, ICML’03*, page 912–919. AAAI Press, 2003.
444 ISBN 1577351894.
- 445 [42] Jason Weston, Frédéric Ratle, and Ronan Collobert. Deep learning via semi-supervised
446 embedding. In *Proceedings of the 25th International Conference on Machine Learning*,
447 pages 1168–1175, New York, NY, USA, 2008. Association for Computing Machinery. ISBN
448 9781605582054.
- 449 [43] Louis-Pascal Xhonneux, Meng Qu, and Jian Tang. Continuous graph neural networks. In
450 *International Conference on Machine Learning*, pages 10432–10441. PMLR, 2020.
- 451 [44] Eda Bayram. Propagation on multi-relational graphs for node regression. In Rosa Maria Benito,
452 Chantal Cherifi, Hocine Cherifi, Esteban Moro, Luis M. Rocha, and Marta Sales-Pardo, editors,
453 *Complex Networks & Their Applications X*, pages 155–167, Cham, 2022. Springer International
454 Publishing. ISBN 978-3-030-93409-5.
- 455 [45] Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Combining neural
456 networks with personalized pagerank for classification on graphs. In *International Conference on*
457 *Learning Representations*, 2019. URL <https://openreview.net/forum?id=H1gL-2A9Ym>.
- 458 [46] Ming Chen, Zhewei Wei, Bolin Ding, Yaliang Li, Ye Yuan, Xiaoyong Du Du, and Ji-Rong Wen.
459 Scalable graph neural networks via bidirectional propagation. 2020.
- 460 [47] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-
461 Rad. Collective classification in network data. *AI Magazine*, 29(3):93–106, 2008.
- 462 [48] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann.
463 Pitfalls of graph neural network evaluation. *Relational Representation Learning Workshop*,
464 *NeurIPS*, 2018.
- 465 [49] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
466 Catasta, and Jure Leskovec. Open Graph Benchmark: Datasets for machine learning on graphs.
467 *arXiv preprint arXiv:2005.00687*, 2020.
- 468 [50] Vijay Prakash Dwivedi, Chaitanya K Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio,
469 and Xavier Bresson. Benchmarking graph neural networks. *arXiv preprint arXiv:2003.00982*,
470 2020.
- 471 [51] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large
472 graphs. In *Proceedings of the 31st International Conference on Neural Information Processing*
473 *Systems*, pages 1025–1035, 2017.
- 474 [52] Johannes Klicpera, Stefan Weissenberger, and Stephan Günnemann. Diffusion improves graph
475 learning. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- 476 [53] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd*
477 *International Conference on Learning Representations, ICLR*, 2015.
- 478 [54] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken ichi Kawarabayashi, and
479 Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In
480 Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on*
481 *Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5453–5462.
482 PMLR, 10–15 Jul 2018.
- 483 [55] Sami Abu-El-Haija, Bryan Perozzi, Amol Kapoor, Hrayr Harutyunyan, Nazanin Alipourfard,
484 Kristina Lerman, Greg Ver Steeg, and Aram Galstyan. Mixhop: Higher-order graph convolution

- 485 architectures via sparsified neighborhood mixing. In *International Conference on Machine*
 486 *Learning (ICML)*, 2019.
- 487 [56] Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Be-
 488 yond homophily in graph neural networks: Current limitations and effective designs. *Advances*
 489 *in Neural Information Processing Systems (NeurIPS)*, 2020.
- 490 [57] Abraham Berman and Robert J. Plemmons. *Nonnegative Matrices in the Mathematical Sciences*.
 491 SIAM, 1994.
- 492 [58] Fan R. K. Chung. *Spectral Graph Theory*. Number 92. American Mathematical Soc., 1997.
- 493 [59] Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric.
 494 In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- 495 [60] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David
 496 Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J.
 497 van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew
 498 R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W.
 499 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.
 500 Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul
 501 van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific
 502 Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.

503 A Appendix

504 A.1 Closed-Form Solution for Harmonic Interpolation

505 Given the *Dirichlet energy* $\ell(\mathbf{x}, G) = \frac{1}{2} \mathbf{x}^\top \Delta \mathbf{x}$, we want to solve for missing features $\mathbf{x}_u =$
 506 $\operatorname{argmin}_{\mathbf{x}_u} \ell$, 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
 507 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
 508 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
 509 (Proposition 3.1). The proof of this result follows from the following Lemma.

510 **Lemma A.1.** *Take any undirected and connected graph with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$, and*
 511 *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*
 512 *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.*

513 *Proof.* Define

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

514 to be the matrix equal to $\tilde{\mathbf{A}}_{uu}$ in the lower right $b \times b$ sub-matrix and padded with zero entries
 515 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
 516 an adjacency matrix of some strongly connected graph and is therefore irreducible [57, Theorem
 517 2.2.7]. These observations allow us to deduce that $\rho(\tilde{\mathbf{A}}_{up}) < \rho(\tilde{\mathbf{A}})$ [57, Corollary 2.1.5]. Note that
 518 $\rho(\tilde{\mathbf{A}}_{up}) = \rho(\tilde{\mathbf{A}}_{uu})$ as $\tilde{\mathbf{A}}_{up}$ and $\tilde{\mathbf{A}}_{uu}$ share the same non-zero eigenvalues. Furthermore, $\rho(\tilde{\mathbf{A}}) \leq 1$
 519 as we can write $\tilde{\mathbf{A}} = \mathbf{I} - \Delta$ and Δ is known to have eigenvalues in the range $[0, 2]$ [58]. Combining
 520 these inequalities gives the result $\rho(\tilde{\mathbf{A}}_{uu}) = \rho(\tilde{\mathbf{A}}_{up}) < \rho(\tilde{\mathbf{A}}) \leq 1$. \square

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

525 *Proof.* To prove Δ_{uu} is non-singular it is enough to show 0 is not an eigenvalue. Note that $\Delta_{uu} =$
 526 $\mathbf{I} - \tilde{\mathbf{A}}_{uu}$ so 0 is not an eigenvalue if and only if $\tilde{\mathbf{A}}_{uu}$ does not have an eigenvalue equal to 1, which
 527 follows from Lemma A.1. \square

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.

528 **A.2 Closed-Form Solution for the Euler scheme**

529 **Proposition A.3.** Take any undirected and connected graph with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$,
 530 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 $\mathbf{x} =$
 531 $\mathbf{x}^{(0)} \in \mathbf{R}^n$ be the initial feature vector and define the following recursive relation

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

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

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

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

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

534 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
 535 convergence of the last u rows

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

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

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

538 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
 539 following limit

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

540 □

541 **A.3 Baselines' Implementation and Tuning**

542 **Label Propagation.** We use the label propagation implementation provided in Pytorch-
 543 Geometric [59]. Since the method is quite sensitive to the value of the α hyperpa-
 544 rameter, we perform a gridsearch separately on each dataset over the following values:
 545 $[0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.99]$.

546 **Positional Encodings.** We compute the laplacian eigenvectors using SciPy [60] sparse eigenvectors
 547 routines. We use the top twenty eigenvectors as positional encodings.

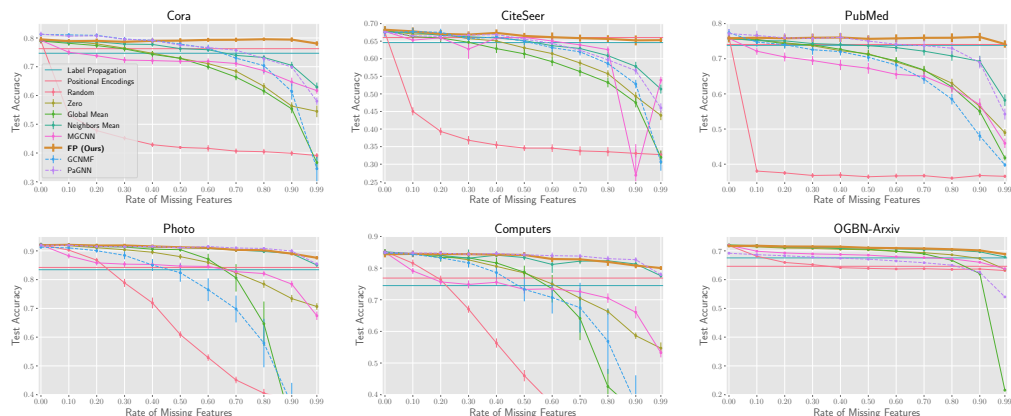


Figure 6: Test accuracy for varying rate of missing features on six common node-classification benchmarks. For methods that require a downstream GNNs, a 2-layer GraphSAGE [51] is used. On OGBN-Arxiv, GCNMF goes out-of-memory and is not reported.

548 **MGCNN.** We re-implement MGCNN [36] in Pytorch by taking inspiration from the authors’ public
 549 TensorFlow code ⁵. For simplicity, we use the version of the model with only graph convolutional
 550 layers and without an LSTM. For the matrix completion training process, we split the observed
 551 features into 50% input data, 40% training targets and 10% validation data. Once the MGCNN
 552 model is trained, we feed it the matrix with all the observed features to predict the whole feature
 553 matrix. This reconstructed features matrix is then used as input for a downstream GNN (as for the
 554 feature-imputation baselines).

555 **A.4 Discussion Over Baselines’ Performance**

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

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

573 However, we find Random to perform better than Zero when all features are missing. This is in line
 574 with findings in the literature [39, 40], where Random features have been shown to work well in
 575 conjunction with GNNs as they act as signatures for the nodes. On the other hand, if all nodes have
 576 all zero vectors, it becomes basically impossible to distinguish them. After applying a GNN, all nodes
 577 will still have very similar embeddings and the task performance will be close to a random guess.

⁵<https://github.com/fmonti/mgcnn>