
Curvature Graph Network

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Graph-structured data is prevalent in many domains. Despite the widely celebrated
2 success of deep neural networks, their power in graph-structured data is yet to
3 be fully explored. We propose a novel network architecture that incorporates
4 advanced graph structural information. In particular, we leverage discrete graph
5 curvature, which measures how the neighborhoods of a pair of nodes are structurally
6 related. The curvature of an edge (x, y) defines the distance taken to travel from
7 neighbors of x to neighbors of y , compared with the length of edge (x, y) . It is a
8 much more descriptive structural measure compared to previously ones that only
9 focus on node specific attributes or limited topological information such as degree.
10 Our curvature graph convolution network outperforms state-of-the-art on various
11 real-world graphs, especially the larger and denser ones.

12 1 Introduction

13 Despite the huge success of deep neural networks, it remains challenging to fully exploit their power
14 on *graph-structured data*, i.e., data whose underlying structure is a graph, e.g., a social network,
15 a telecommunication network and a biological network. Inspired by the power of convolution on
16 image data, there are many attempts to extend convolutional networks for graph-structured data. For
17 example, one category of graph neural networks, named spatial approaches, iteratively update each
18 node representation by aggregating the information from its neighbors. For these spatial approaches,
19 it is important to incorporate local structural information of the graph. Node degree has been used
20 to reparametrize the nonlinear transformation of messages [9] or as an additional node feature [4].
21 However, node degree is fairly limited; there can be different graph topologies with the same degree
22 distribution. The limitation is illustrated in Figure 1. Nodes x and y have the same degree in three
23 significantly different graphs: a tree, a grid graph and a clique. You *et al.* [20] calculates the shortest
24 path between nodes. However this is not for node classification task and cannot extract complicate
25 local information. To effectively make use of graph structural knowledge, one would need a feature
26 with more discriminative power; one that can distinguish these three scenarios in Figure 1.

27 In this paper, we are focusing on designing a novel graph neural network that exploits advanced
28 structural information. Notice that node degree only describes the number of neighbors of each node,
29 but does not say how these neighbors are connected among themselves. We seek to use structural
30 information characterizing how neighborhoods of a pair of nodes relate to each other. In Figure 1, the
31 neighborhoods of x and y are well separated in a tree. In a grid graph, the two neighborhoods are
32 within a parallel shift of each other. In a clique, they completely overlap. To quantify such pairwise
33 structural information, we draw inspiration from the recent study of *graph curvature* [12, 7, 18].

34 Similar to the curvature in the continuous domain, e.g., the Ricci curvature of a Riemannian manifold,
35 the discrete graph curvature measures how the geometry of a pair of neighborhoods deviates from a
36 “flat” case, namely, the case of a grid graph. There are several definitions of discrete curvature for
37 graphs. The most notable ones are Ollivier’s Ricci curvature [12] and Forman curvature [18]. In both
38 definitions, the edges of a (infinite) grid graph have zero curvature. The curvature of an edge (x, y)
39 in a tree is negative and edges in a complete graph have positive curvature. Intuitively, the graph

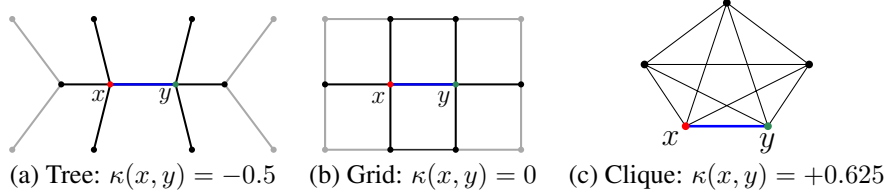


Figure 1: Illustration of structural information. In all three graphs, the degrees of x and y are the same. However, the Ricci curvature of the edge (x, y) is negative, zero, and positive, respectively. All edges have weight 1.

40 curvature measures how well two neighborhoods are connected and/or overlap with each other. Such
 41 information is related to how information propagates in the neighborhood, and should be leveraged
 42 by a graph convolutional network.

43 We propose *Curvature Graph Network (CurvGN)*, which is built on advance graph curvature informa-
 44 tion. In particular, the proposed CurvGN is a novel network architecture that efficiently computes
 45 graph curvature and fully leverages such information in graph convolution. Using curvature informa-
 46 tion, CurvGN better adapts to different local structural scenarios and filter messages passed between
 47 nodes differently. Notice that the curvature information captures how easy information flows between
 48 the nodes. Within a well-connected community, neighborhoods of adjacent nodes have large overlap
 49 and many shortcuts. The corresponding curvature is positive and passing information between the
 50 nodes is easy. For edges bridging two clusters/cliques, the curvature is negative and information
 51 is hard to pass. A key to our success is that we choose to be agnostic on whether the curvature
 52 information should be used to block or accelerate the messages in graph convolution. We exploit
 53 the curvature in a data-driven manner and learn how to use it to reweigh different channels of the
 54 message.

55 2 Curvature Graph Network

56 We first formulate the node label prediction problem of a graph, and explain the mechanism of a
 57 Graph Neural Network (GNN). Suppose we have an undirected graph $G = (V, E)$ with features on
 58 the vertices $H = (h_1, h_2, \dots, h_n), h_i \in \mathbb{R}^F$. Here $n = |V|$ is the number of nodes in the graph
 59 and F is the feature dimension of each node. Given labels of some nodes in V , we would like to
 60 predict the labels of the remaining nodes. A GNN consists of multiple hidden layers that update
 61 node representation from lower level node representation $H^t \in \mathbb{R}^{n \times F_t}$ to high level representation
 62 $H^{t+1} \in \mathbb{R}^{n \times F_{t+1}}$. In particular, H^0 is the input feature, H . Node representations of the last layer,
 63 H^T , are fed to a fully connected layer or a linear classifier to predict node labels. The layers and
 64 their representations are illustrated in the top of Figure 2.

65 Now we explain how to construct hidden layers that update node representations from H^t to H^{t+1} .
 66 We focus on spatial approaches and treat the convolution as a message passing scheme. The $(t+1)$ -th
 67 representation of node x is computed by aggregating messages passed from x 's neighbors. We also
 68 include the message from x to itself. There are several aggregation methods, such as mean, max and
 69 sum. We choose summation as it is a commonly used aggregation method [5, 16, 19]. Denote by
 70 $\bar{\mathcal{N}}(x) = \mathcal{N}(x) \cup \{x\}$ the neighborhood of x including itself. We have $h_x^{t+1} = \sigma_t \left(\sum_{y \in \bar{\mathcal{N}}(x)} M_{y \rightarrow x}^t \right)$,
 71 in which σ_t is the non-linear transformation. A message passed from y to x is a linear transformation
 72 of y 's representation. We also introduce a weight τ_{xy}^t whose purpose will be clear later. We have
 73 $M_{y \rightarrow x}^t = \tau_{xy}^t W^t h_y^t$, in which W^t is the linear transformation matrix learned in training. Formally,
 74 we have the representation updating equation

$$h_x^{t+1} = \sigma_t \left(\tau_{xx}^t W^t h_x^t + \sum_{y \in \mathcal{N}(x)} \tau_{xy}^t W^t h_y^t \right) \quad (2.1)$$

75 It is crucial to obtain suitable reweighting parameter τ_{xy}^t since it is directly affecting how neighboring
 76 node information are passed to the node x . Some papers use node degree information as τ_{xy}^t [5, 9]
 77 and other work uses joint node features to compute the self attention as τ_{xy}^t [16]. We propose to
 78 use more advanced structural information, i.e., the Ricci curvature, to compute τ_{xy}^t . It is also worth
 79 mentioning that the reweighting parameter τ_{xy}^t is not necessarily a scalar. It can also be anything
 80 between a scalar and a $F^t \times F^t$ matrix. In fact, we choose F^t vector later on because it has more
 81 expressive power than a scalar and it is easier to train than a matrix.

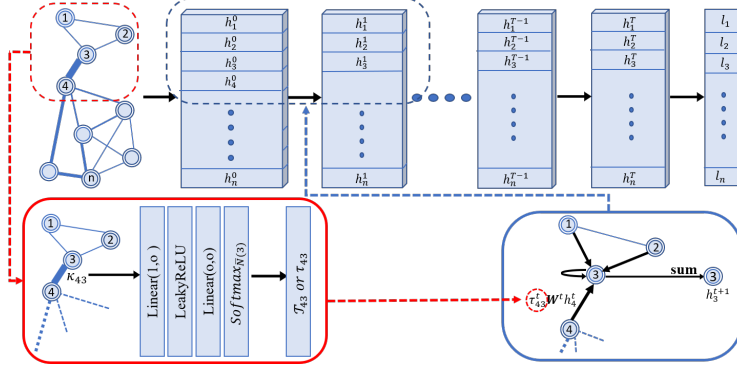


Figure 2: An overview of our Curvature Graph Network.

82 To illustrate how we build curvature convolution layer in Equation (2.1), we define Ricci Curvature in
 83 the context of graph (Section 2.1). We explain how to compute τ_{xy}^t from the curvature in Section 2.2.

84 2.1 Ricci Curvature

85 In Riemannian geometry, curvature measures how a smooth object deviates from being flat, or being
 86 straight in the case of a line. Similar concepts can be extended to non-smooth setting for discrete
 87 objects. In particular, curvature has been studied for metric-measure space in [1, 2, 8, 15], Markov
 88 chain by Ollivier [12] and general graphs in [7]. The definitions of curvatures that are easier to
 89 generalize in a discrete graph setting are sectional curvature and Ricci curvature.

90 To generalize Ricci curvature to discrete spaces, Ollivier [12] takes a coarse approach that represents
 91 S_x as a probability measure m_x of mass 1 around x . Thus the distance can be measured by Wasser-
 92 stein distance (or Earth Mover distance) which finds the optimal mass-preserving transportation
 93 plan between two probability measures. Then the coarse Ricci curvature $\kappa(x, y)$ on edge (x, y)
 94 is defined by comparing the Wasserstein distance $W(m_x, m_y)$ to the distance $d(x, y)$, formally,
 95 $\kappa_{xy} = 1 - W(m_x, m_y)/d(x, y)$. For an undirected graph $G = (V, E)$, denotes the set of neighbor-
 96 ing nodes of a node $x \in V$ as $N(x) = \{x_1, x_2, \dots, x_k\}$. Then we can define a probability measure
 97 m_x^α at x : $m_x^\alpha(x_i) = \delta_{x_i=x} \times \alpha + \delta_{x_i \in N(x)} \times (1 - \alpha)/k$ where α is a parameter within $[0, 1]$ and δ
 98 is the indicator function. It is to keep probability mass of α at node x itself and distribute the rest
 99 uniformly over the neighborhood. To compute the Wasserstein distance $W(m_x^\alpha, m_y^\alpha)$ between the
 100 probability measures around two end points x, y of the edge (x, y) , the optimal transportation plan
 101 can be solved by the following linear programming:

$$\begin{aligned} & \min_M \sum_{i,j} d(x_i, y_j) M(x_i, y_j), & (2.2) \\ & \text{s.t.} \sum_j \forall i, M(x_i, y_j) = m_x^\alpha(x_i); \sum_i \forall j, M(x_i, y_j) = m_y^\alpha(y_j). \end{aligned}$$

102 where $M(x_i, y_j)$ is the amount of probability mass transported from node x_i to y_j along the shortest
 103 path with length $d(x_i, y_j)$. Following existing work [11], we set $\alpha = 0.5$.

104 2.2 Curvature-Driven Graph Convolution

105 Next we present how Ricci curvature is used in our graph convolutional network. The usage of
 106 curvature should depend on the problem and the data. Intuitively, curvature measures how easy a
 107 message flows through an edge, and should be used to control messages in convolution. For example,
 108 an edge with negative curvature is likely to be a bridge connecting two different communities. If
 109 we assume different communities tend to have different representations/labels, a message should be
 110 blocked on this edge. Meanwhile, an edge with positive curvature tends to be intra-community and
 111 thus should have accelerated message flow. However, the intuition may be invalid if the community
 112 structure is not correlated with node representation/labels.

113 We choose to be agnostic on how the knowledge of edge curvature should be used. We resort to
 114 a data-driven strategy and learn a mapping function that maps Ricci curvature κ_{xy} to the weight
 115 of messages, i.e., τ_{xy}^t in Equation (2.1). We first explain how the mapping is learned end-to-end
 116 (CurvGN-1). Next we expand the mapping to a multi-valued version, to incorporate more flexibility
 117 in the model (CurvGN-n).

118 **CurvGN-1.** As mentioned before, τ_{xy}^t can be anything between a scalar and a $F^t \times F^t$ matrix. We
 119 first assume τ_{xy}^t is a scalar. In this case, the mapping function can be defined as: $f^t : \kappa_{xy} \rightarrow \tau_{xy}^t$.

120 We create a multi-layer perceptron (MLP) to approximate the mapping function f^t since MLP is
 121 proved to be a universal approximation machine and can be easily incorporated into our GNN model
 122 for end-to-end training. Denote the MLP at the t -th layer as MLP^t . As summation is used as the
 123 aggregation function in Equation (2.1), the messages may accumulate to an arbitrarily large value. To
 124 prevent a numerical explosion, we use softmax to normalize outputs of MLP over all neighbors. This
 125 gives us the eventual weight, τ_{xy}^t . Figure 2 bottom shows how the MLP transforms a curvature and
 126 uses it to reweigh messages.

127 **CurvGN-n.** Messages $M_{y \rightarrow x}^t$ are usually multi-channeled. In particular, they are F^{t+1} -dimensional.
 128 The scalar weight generated using curvature is not necessarily the same for different channels. To
 129 improve the expressing power of τ_{xy}^t , we create a similar mapping function as f^t . But the new
 130 mapping generates a reweighing vector $\mathcal{T}_{xy}^t \in \mathbb{R}^{F^{t+1}}$. In other words, we learn to reweigh different
 131 message channels differently. More details can be found in the appendix.

132 3 Experiments

133 Our real-world benchmarks include two families of datasets: small sparse graphs and large dense
 134 graphs. We compare our networks CurvGN-1 and CurvGN-n with several strong baselines. Aside
 135 from commonly compared baselines GCN and GAT, we also compare CurvGN-1 and CurvGN-n with
 136 multilayer perceptron (MLP), MoNet [9] and GraphSAGE with mean aggregation (GS-mean) [4].
 137 Our method is on par with state-of-the-art methods on relatively small graphs and greatly outperforms
 138 state-of-the-art methods on large and dense graphs, which tend to have heterogeneous topology.

139 **Datasets.** We use three popular citation network benchmark datasets: Cora, Citeseer and PubMed
 140 [13]. We categorize Cora and Citeseer into the first family since both Cora and Citeseer graphs are
 141 relatively small and sparse. They have thousands of nodes and edges with an average node degree
 142 below 2. We also use four extra datasets in [14]: Coauthor CS, Coauthor Physics, Amazon Computers
 143 and Amazon Photos. These graphs, together with PubMed, are large and dense graphs. Those graphs
 144 have more than 10 thousands node and 200 thousands edges with an average node degree as high as
 145 20. Descriptions and statistics for all datasets in our experiments can be found in the Appendix.

146 In Table 1, we report the mean and standard deviation of classification accuracy on test nodes on 100
 147 runs and re-use the metrics reported by [9, 14, 16] for other state-of-the-art methods.

Table 1: Performance on Real-World Benchmarks

Method	Cora	Citeseer	PubMed	Coauthor CS	Coauthor Physics	Amazon Computer	Amazon Photo
MLP	58.2	59.1	70.0±2.1	88.3±0.7	88.9±1.1	44.9±5.8	69.6±3.8
MoNet	81.7	71.2	78.6±2.3	90.8±0.6	92.5±0.9	83.5±2.2	91.2±1.3
GS-mean	79.2	71.2	77.4±2.2	91.3±2.8	93.0±0.8	82.4±1.8	91.4±1.3
GCN	81.5±0.5	70.9±0.5	79.0±0.3	91.1±0.5	92.8±1.0	82.6±2.4	91.2±1.2
GAT	83.0±0.7	72.5±0.7	79.0±0.3	90.5±0.6	92.5±0.9	78.0±19.0	85.1±20.3
CurvGN-1	82.6±0.6	71.5±0.8	78.8±0.6	92.9±0.4	94.1±0.3	86.3±0.7	92.5±0.5
CurvGN-n	82.7±0.7	72.1±0.6	79.2±0.5	92.8±0.3	94.3±0.2	86.5±0.7	92.5±0.5

148 **Discussion.** Our method is on par with state-of-the-art performance for relatively small graphs.
 149 Meanwhile, it achieves superior performance on large and dense graphs. It is clear curvature
 150 information provides richer information than node degrees, at least on large and dense graphs.
 151 We also observe that CurvGN-n is generally better than CurvGN-1. This means multi-channel
 152 reweighting provides a better mechanism in leveraging curvature information. These observations
 153 are consistent with our synthetic experiments, which are not included due to page limitation. We
 154 conduct experiments on synthetic datasets generated according to various well-established graph
 155 models, e.g., stochastic block model [3], Watts-Strogatz network [17], Newman-Watts network [10]
 156 and Kleinberg’s navigable small world graph [6]. On these data, CurvGN consistently outperforms
 157 GAT and GCN, demonstrating the benefit of curvature information.

References

- 158
- 159 [1] Anca-Iuliana Bonciocat. A rough curvature-dimension condition for metric measure spaces.
160 *Open Mathematics*, 12(2):362–380, 2014.
- 161 [2] Anca-Iuliana Bonciocat and Karl-Theodor Sturm. Mass transportation and rough curvature
162 bounds for discrete spaces. *Journal of Functional Analysis*, 256(9):2944–2966, 2009.
- 163 [3] Aurelien Decelle, Florent Krzakala, Cristopher Moore, and Lenka Zdeborová. Asymptotic
164 analysis of the stochastic block model for modular networks and its algorithmic applications.
165 *Physical Review E*, 84(6):066106, 2011.
- 166 [4] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
167 graphs. In *Advances in Neural Information Processing Systems*, pages 1024–1034, 2017.
- 168 [5] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
169 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 170 [6] Jon M Kleinberg. Navigation in a small world. *Nature*, 406(6798):845, 2000.
- 171 [7] Yong Lin, Linyuan Lu, and Shing-Tung Yau. Ricci curvature of graphs. *Tohoku Mathematical
172 Journal, Second Series*, 63(4):605–627, 2011.
- 173 [8] John Lott and Cédric Villani. Ricci curvature for metric-measure spaces via optimal transport.
174 *Annals of Mathematics*, pages 903–991, 2009.
- 175 [9] Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and
176 Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model
177 cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
178 pages 5115–5124, 2017.
- 179 [10] Mark EJ Newman and Duncan J Watts. Renormalization group analysis of the small-world
180 network model. *Physics Letters A*, 263(4-6):341–346, 1999.
- 181 [11] Chien-Chun Ni, Yu-Yao Lin, Jie Gao, and Xianfeng Gu. Network alignment by discrete ollivier-
182 ricci flow. In *International Symposium on Graph Drawing and Network Visualization*, pages
183 447–462. Springer, 2018.
- 184 [12] Yann Ollivier. Ricci curvature of markov chains on metric spaces. *Journal of Functional
185 Analysis*, 256(3):810–864, 2009.
- 186 [13] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-
187 Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- 188 [14] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann.
189 Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.
- 190 [15] Karl-Theodor Sturm et al. On the geometry of metric measure spaces. *Acta mathematica*,
191 196(1):65–131, 2006.
- 192 [16] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
193 Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- 194 [17] Duncan J Watts and Steven H Strogatz. Collective dynamics of ‘small-world’ networks. *nature*,
195 393(6684):440, 1998.
- 196 [18] Melanie Weber, Jürgen Jost, and Emil Saucan. Forman-ricci flow for change detection in large
197 dynamic data sets. *Axioms*, 5(4):26, 2016.
- 198 [19] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
199 networks? *arXiv preprint arXiv:1810.00826*, 2018.
- 200 [20] Jiaxuan You, Rex Ying, and Jure Leskovec. Position-aware graph neural networks. *arXiv
201 preprint arXiv:1906.04817*, 2019.

202 **A Appendix**

203 **Details of CurvGN-n.** Using the same strategy as CurvGN-1, the vector \mathcal{T}_{xy}^t is calculated by
 204 applying a MLP^t with F^{t+1} outputs. Then, we apply a channel-wise softmax function, \mathbf{S}^t , that
 205 normalizes the MLP outputs separately on each message channel: $\mathcal{T}_{xy}^t = \mathbf{S}^t(\text{MLP}^t(\kappa_{xy}))$

206 Substituting \mathcal{T}_{xy} into Equation (2.1), we have the convolution of CurvGN-n:

$$h_x^{t+1} = \sigma_t \left(\sum_{y \in \bar{\mathcal{N}}(x)} \text{diag}(\mathcal{T}_{xy}^t) W^t h_y^t \right) \tag{A.1}$$

207 Here $\text{diag}(\mathcal{T}_{xy}^t)$ is a matrix whose diagonal entries are entries of \mathcal{T}_{xy}^t .

208 **Design details of the network.** In practice, we use a two-convolutional-layer CurvGN model. The
 209 first layer is a linear transform layer that produces an output feature vector paired with a three
 210 layer MLP that computes reweighing vector. input layer that transforms curvature information into
 211 64 dimension feature vector and a linear layer followed by a LeakyReLU layer that generates a
 212 reweighing vector of dimension 64. The output feature is pushed into an exponential linear unit layer
 213 to add non-linearity. The second layer is for classification, with the same structure as the first layer
 214 except that the output feature is now at length of class number. The hyperparameters are similar to
 215 GAT implemented in [16].

216 **Statistical detail of benchmarks** We describe the statistical details of all datasets in Table 2. Cora
 217 and Citeseer are considered as small and sparse graphs while PubMed, Coauthors and Amazons are
 218 considered as large and dense graphs.

Table 2: Statistic details of all datasets.

Datasets	#Classes	#Nodes	#Edges	#Features	#Training	#Edges/#Nodes
Cora	7	2708	5429	1433	140	2.0
Citeseer	6	3327	4732	3703	120	1.42
PubMed	3	19717	44338	500	60	2.25
Coauthor CS	15	18333	100227	6805	300	5.47
Coauthor Physics	5	34493	282455	8415	100	8.19
Amazon Computers	10	13381	259159	767	200	19.37
Amazon Photo	8	7487	126530	745	160	16.90

219 **Data splitting and Hyper-parameters** We use the exact data splitting as in semi-supervised learning
 220 setting used in [5, 16]: using 20 nodes per class for training, 500 nodes for validation and 1000 nodes
 221 for testing. During training stage, we set L_2 regularization with $\lambda = 0.0005$ for all datasets. All the
 222 models are initialized by Glorot initialization and trained by minimizing cross-entropy loss using
 223 Adam SGD optimizer with learning rate $r = 0.005$. We apply an early stopping strategy with the
 224 help of the validation set based on the validation set’s accuracy with a patience of 100 epochs.