# Node Copying for Protection Against Graph Neural Network Topology Attacks

Florence Regol, Soumyasundar Pal and Mark Coates

Department of Electrical and Computer Engineering

McGill University, 3480 University Street, Montreal, Quebec, Canada H3A2A7

Email: {florence.robert-regol, soumyasundar.pal}@mail.mcgill.ca, mark.coates@mcgill.ca

*Abstract*—**Adversarial attacks can affect the performance of existing deep learning models. With the increased interest in graph based machine learning techniques, there have been investigations which suggest that these models are also vulnerable to attacks. In particular, corruptions of the graph topology can degrade the performance of graph based learning algorithms severely. This is due to the fact that the prediction capability of these algorithms relies mostly on the similarity structure imposed by the graph connectivity. Therefore, detecting the location of the corruption and correcting the induced errors becomes crucial. There has been some recent work which tackles the detection problem, however these methods do not address the effect of the attack on the downstream learning task. In this work, we propose an algorithm that uses node copying to mitigate the degradation in classification that is caused by adversarial attacks. The proposed methodology is applied only after the model for the downstream task is trained and the added computation cost scales well for large graphs. Experimental results show the effectiveness of our approach for several real world datasets.**

*Index Terms*— Adversarial attacks, Graph convolutional networks, Semi-supervised learning

## I. INTRODUCTION

The application of deep learning models in real world systems has become increasingly prevalent, and as a result there has been an increased attention paid to their robustness and vulnerability to adversarial attack [1]. It has been demonstrated that many deep neural networks are susceptible to malicious attack and this has given rise to serious concerns regarding their reliability.

In many problem domains, including recommender systems, fraud detection, disease outcome and drug interaction prediction, there are structural relationships between data items. A graph is a natural mechanism for representing these relationships and this has led to the desire to translate the success of neural networks to the graph setting. An intense research effort has led to many models and algorithms [2]–[9]. It has been demonstrated that knowledge of the graph can be leveraged to compensate for having limited access to labelled data. Subsequently, there has been successful industrial application of these models [10]–[12]. This has raised concerns regarding the vulnerabilities of graph neural networks (GNNs) and researchers have commenced the development and investigation of attacks and defence mechanisms. Understanding the adversarial vulnerabilities of GNNs helps to expose the limitations of existing GNN models and can inspire better models and training strategies [13], [14].

Convolutional neural networks are usually subject to attacks that involve data manipulation to alter features. Graph neural networks can be targeted by similar attacks, but they are also subject to an alternative form of attack that involves alteration of the graph topology. In [15], Zügner et al. proposed Nettack, a method for constructing adversarial perturbations of graph data, which alters the graph topology and/or the node attributes in order to produce significant degradation in node classification performance. The experimental analysis in [15] suggests that attacks on the graph topology can have a more severe impact on classification performance compared to feature alteration. The attack in [15] strives to disrupt the classification of individual nodes in the graph; more recent work has targeted the deterioration of performance across the entire graph [16]. Other adversarial attacks on graphs have been proposed that highlight the vulnerability of GNNs for a wider range of inference tasks. In [17], Dai et al. show the efficacy of their proposed method on a real-world financial dataset where the classification task is to distinguish normal transactions from abnormal ones. This practical scenario gives a concrete example of how harmful such attacks can be and motivates the need for designing efficient countermeasures.

In response to the development of attacks on graph learning, there has been some preliminary research into detecting attacks. Zhang et al. [18] propose an algorithm to detect which nodes have been subjected to an attack via modification of their edges. The procedure relies on the inconsistencies that the attack induces in the classification outputs in the neighbourhood of an attacked node. Although the technique in [18] offers a promising (albeit not foolproof) approach for detection of an attack, it does not provide a mechanism for rectifying the output of the learning algorithm.

In this paper, we focus on the next stage in the learning pipeline. We address the question of what to do after a detection procedure has notified us that there is a high probability that a node has been subjected to a topology attack. We introduce a copying procedure to partially recover the model accuracy of a graph convolutional neural network (GCN) for the corrupted nodes. The procedure involves copying the features of an attacked node to multiple similar locations in the graph and evaluating the output at these locations. The intuition is that when the features are moved to locations that correspond to its true class and have not been attacked, the GCN will return a correct classification. Through analysis of

citation network datasets, we illustrate that this procedure can improve the classification accuracy by 10-15 percent for the attacked nodes.

The paper is organized as follows. In Section II, we present background material, briefly reviewing graph convolutional neural networks. Section III provides more detail regarding the problem setting and Section IV presents our proposed recovery methodology. Section V describes the numerical experiments and presents and discusses the results. Section VI concludes the paper and suggests future research directions.

## II. GRAPH CONVOLUTIONAL NETWORKS

For the scope of this paper, we address a downstream node classification task using a GCN proposed in [2], [19]. In this setting, we are given a set of nodes $\mathcal{V}$ and edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ that form a graph $\mathcal{G}_{obs} = (\mathcal{V}, \mathcal{E})$. Node $i$ is associated with a feature vector $\boldsymbol{x}_i$ and a label $\boldsymbol{y}_i$.

In the semi-supervised setting, we have knowledge of labels only at a limited subset of nodes, $\mathcal{V}_{train} \subset \mathcal{V}$, and we aim to predict the labels at the nodes at the test set, $\mathcal{V}_{test} \subset \mathcal{V}$. The model uses information provided by the observed graph $\mathcal{G}_{obs}$, the complete feature matrix $\mathbf{X} = [\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N]^T$ and the labels in the training set $\mathbf{Y}_{\mathcal{V}_{\mathbf{train}}} = \{\boldsymbol{y}_i : i \in \mathcal{V}_{train}\}$.

The layerwise propagation rule in simpler GCN architectures [2], [19] is based on a graph convolution operation and can be written as:

$$\mathbf{H}^{(1)} = \sigma(\hat{\mathbf{A}}_\mathcal{G} \mathbf{X} \mathbf{W}^{(0)}), \tag{1}$$

$$\mathbf{H}^{(l+1)} = \sigma(\hat{\mathbf{A}}_\mathcal{G} \mathbf{H}^{(l)} \mathbf{W}^{(l)}). \tag{2}$$

Multiplication with the normalized adjacency operator $\hat{\mathbf{A}}_\mathcal{G}$ results in aggregation of the output features across node neighborhood at each layer. $\mathbf{W}^{(l)}$ is a matrix of trainable weights at layer $l$ of the neural network and $\sigma$ denotes a pointwise non-linear activation function. $\mathbf{H}^{(l)}$ is the output representation from layer $l - 1$. In a node classification setting, using a $L$-layer network, the prediction is obtained by applying a softmax activation in the last layer and is written as $\hat{\mathbf{Y}} = \mathbf{H}^{(L)}$. The weights of the neural network are learned via backpropagation with the objective of minimizing the cross entropy loss between the training labels $\mathbf{Y}_{\mathcal{V}_{train}}$ and the network predictions $\hat{\mathbf{Y}}_{\mathcal{V}_{train}} = \{\hat{\boldsymbol{y}}_i : i \in \mathcal{V}_{train}\}$ at the nodes in the training set.

## III. PROBLEM SETTING

As stated in Section II, we address a semi-supervised node classification task. Based on the graph $\mathcal{G}_{obs}$, node features $\mathbf{X}$ and a small subset of known training labels $\mathbf{Y}_{\mathcal{V}_{train}}$, the goal is to infer the labels of the nodes in the test set, $\mathcal{V}_{test} = \mathcal{V} \setminus \mathcal{V}_{train}$. However, a subset of the test nodes $\mathcal{V}_{attacked} \subset \mathcal{V}_{test}$ are subjected to adversarial attack (details in Section III-A), which modifies the graph. We consider a random poisoning attack [15] scenario where the attack precedes the model training. As a result, we only have access to the attacked graph $\mathcal{G}_{attacked}$. We assume that the attack only targets a small number of nodes compared to the size of the whole graph and it does not affect any nodes in $\mathcal{V}_{train}$.

For the scope of this work, we also assume that the identities of the nodes in $\mathcal{V}_{attacked}$ are known. In practice, this does not impose any serious restriction on the applicability of the proposed methodology since any reasonably accurate detection algorithm, such as the one proposed in [18] can be employed to identify the nodes in $\mathcal{V}_{attacked}$. If a node is incorrectly labelled as attacked, our proposed procedure in most cases does not modify the classification output. Our goal is to correct the possible classification errors for the nodes in $\mathcal{V}_{attacked}$ after the poisoning attack has occurred.

### A. DICE Attack

Since the impressive performance of most graph based learning algorithms stems from the presence of edges between similar nodes, we consider an attack which aims to disrupt the similarity structure imposed by the graph connectivity. The DICE (Delete Internally Connect Externally) attack is a simple yet effective random attack. It is parameterized by $0 < \beta \leq 1$, which dictates the severity of the degradation of the nodes in $\mathcal{V}_{attacked}$. We assume that the attacker has complete knowledge of the true labels of the nodes in $\mathcal{G}_{obs}$. For each attacked node $v$ with degree $d_v$, the attacker removes $\lceil \beta d_v \rceil$ of its existing edges at random and inserts new edges between node $v$ and $\lceil \beta d_v \rceil$ other nodes, sampled uniformly from the set of all nodes with different true labels from $v$. As a result, node $v$ has at most $\lfloor (1 - \beta) d_v \rfloor$ neighbours with the same label after the attack. Since this attack does not perturb the degree of the target node, the degree distribution of the nodes in $\mathcal{V}_{attacked}$ remains unaltered.

## IV. METHODOLOGY

In our correction strategy, the classification of each node in $\mathcal{V}_{attacked}$ is performed in the following way. First we train a base GCN classifier using the small subset of labeled nodes $\mathbf{Y}_{\mathcal{V}_{train}}$ on $\mathcal{G}_{attacked}$ and store the obtained model. Then we compute a lower dimension representation of the nodes of $\mathcal{G}_{attacked}$ using a node embedding algorithm. This procedure summarizes the information provided by the graph connectivity and the node features in the embedding. In our experiments, we use the Graph Variational Auto-Encoder (GVAE) [3] to obtain the embeddings but any other suitable techniques can also be employed.

We form a symmetric, pairwise distance matrix $D \in \mathbf{R}_+^{N \times N}$, whose $(i, j)$-th entry is defined as:

$$D_{i.j} = \|e_i - e_j\|_2. \tag{3}$$

Here $e_i$ is the embedding of node $i$ and $D_{i,j}$ is the distance between node $i$ and $j$. This distance matrix $D$ is subsequently used to select a set of similar nodes for each node in $\mathcal{V}_{attacked}$. For node $v$, we form the set $\mathcal{V}_v^{(p)}$ of the $p$ most similar nodes based on the $p$ lowest distances from node $v$ as follows:

$$\mathcal{V}_v^{(p)} = \{1 \leq k \leq N \mid k \neq v, D_{i,k} = D_{i,(j)}, 1 < j \leq p + 1\}, \tag{4}$$

where $D_{i,(j)}$ is the $j$-th order statistic of $\{D_{i,k}\}_{k=1}^N$. Then for each node $k \in \mathcal{V}_v^{(p)}$, we copy the feature at node $v$ to

node $k$ (which is equivalent to copying the $v$-th row of $\mathbf{X}$ to the $k$-th row) and compute the prediction $\hat{\boldsymbol{y}}_{v \to k}$ at node $k$ using the existing GCN model. The prediction for node $v$ using the proposed copying procedure is obtained by computing the average of $\{\hat{\boldsymbol{y}}_{v \to k}\}_{k \in \mathcal{V}_v^{(p)}}$. Figure 1 presents an overview of the complete procedure.
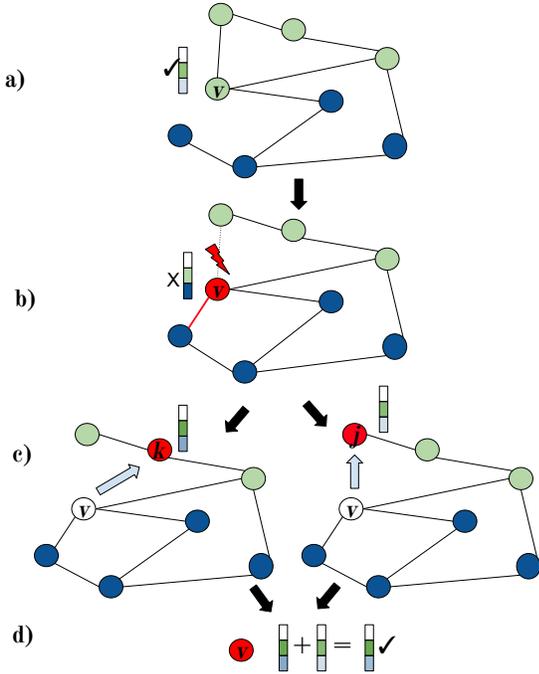


Fig. 1: Summary of the node copying procedure. **a)** In the absence of the attack, the softmax of node $v$ achieves the correct classification in $\mathcal{G}_{obs}$. **b)** Node $v$ is targeted by an attack and is now wrongly classified. **c)** The feature of node $v$ is copied to two new positions $k$ and $j$ and the softmax at those positions $\hat{\boldsymbol{y}}_{v \to k}$ and $\hat{\boldsymbol{y}}_{v \to j}$ are obtained. **d)** The error on node $v$ is corrected by computing the average of $\hat{\boldsymbol{y}}_{v \to k}$ and $\hat{\boldsymbol{y}}_{v \to j}$.

This correction procedure is successful if, on average, the true class of node $v$ is dominating in the set of softmax outcomes. The intuition is that some of the similar nodes included in $\mathcal{V}_v^{(p)}$ will have the same class as the true class $t$ of node $v$. We take the simplified view that, if uncorrupted, a node will often have the same class as most of its neighbors. If the node $v$ is copied at node $k$ in a "wrong" neighborhood, meaning that node $v$ has a different class $w \neq t$ from node $k$ and its neighbors, then pulling the classification of node $v$ to the wrong class $w$ is usually harder and the $w$-th entry in the resulting softmax $\hat{\boldsymbol{y}}_{v \to k}$ is likely to be smaller. However when a node is placed in a "good" neighborhood, the weighted average operation should reinforce the model confidence in the correct class $t$. So when we perform the average over softmax outputs at similar nodes, correct classification can be recovered.

A naive implementation of this method requires $p|\mathcal{V}_{attacked}|$ additional GCN evaluations after the training. This computational burden might be prohibitive for large graphs. However, we note that the prediction at any particular node from an $L$-

layer GCN is influenced only by the $L$-hop neighbourhood of the node, which allows a much cheaper, localized computation of $\{\hat{\boldsymbol{y}}_{v \to k}\}$. The procedure is summarized in Algorithm 1.

---

**Algorithm 1** Error correction using node copying

---

1: **Input:** $\mathcal{G}_{attacked}$, $\mathbf{X}$, $\mathbf{Y}_{\mathcal{V}_{\mathbf{train}}}$, $\mathcal{V}_{attacked}$
2: **Output:** $\widehat{\mathbf{Y}}_{attacked}^{GCN}$, $\widehat{\mathbf{Y}}_{attacked}^{Copying}$
3: Train a GCN using $\mathcal{G}_{attacked}$, $\mathbf{X}$, $\mathbf{Y}_{\mathcal{V}_{\mathbf{train}}}$ and compute $\widehat{\mathbf{Y}}_{attacked}^{GCN} = \{\hat{\boldsymbol{y}}_v\}_{v \in \mathcal{V}_{attacked}}$.
4: Train a GVAE to obtain node embeddings, $\{e_i\}_{i=1}^N$. Compute the pairwise distance matrix $D$ using eq. (3).
5: **for** $v \in \mathcal{V}_{attacked}$ **do**
6:     Form the set $\mathcal{V}_v^{(p)}$ using eq. (4)
7:     **for** $k \in \mathcal{V}_v^{(p)}$ **do**
8:         Copy the features of node $v$ in place of node $k$ and compute $\hat{\boldsymbol{y}}_{v \to k}$ using the existing GCN weights.
9:     **end for**
10:     Compute $\hat{\boldsymbol{y}}_v^{Copying} = \frac{1}{p} \sum_{k \in \mathcal{V}_v^{(p)}} \hat{\boldsymbol{y}}_{v \to k}$
11: **end for**
12: Form $\widehat{\mathbf{Y}}_{attacked}^{Copying} = \{\hat{\boldsymbol{y}}_v^{Copying}\}_{v \in \mathcal{V}_{attacked}}$

---

## V. EXPERIMENTS

We conduct experiments on three citation datasets: Pubmed, Citeseer and Cora [20]. The prediction task is to classify the topics of research articles. Each document is represented as a node in a graph that is formed by adding an edge between any two articles if one of them cites the other. The features consist of a bag-of-words vectors extracted from the contents of the articles. Statistics of the datasets can be seen in Table I.

The purpose of our attack correction algorithm is to retain the advantages derived from the model's ability to exploit knowledge of the graph topology. The information from the graph is more valuable when the amount of labelled data is severely limited, so we focus on this setting.

For each trial, $\mathcal{V}_{train}$ is formed by randomly sampling 10 or 20 nodes per class. Then an additional 50 nodes are sampled from the remaining set of nodes and these are targeted by the attack. The rows in the adjacency matrix of $\mathcal{G}_{obs}$ of each node in $\mathcal{V}_{train}$ are iteratively corrupted following the DICE attack described previously. We consider the parameters $\beta = 0.5$ and $\beta = 0.75$ for the attack to test the robustness of the recovery procedure. The number of new positions for a node $p$ is set to 10 in our experiment. This parameter can be chosen more judiciously through cross-validation. The GCN and GVAE hyperparameters are set to the values specified in [3] and [19], respectively. These are obtained by optimizing the classification accuracy on a validation set of 500 nodes on the Cora dataset.

TABLE I: Datasets statistics

| Dataset | Nodes | Classes | Edges | Features |
|---|---|---|---|---|
| Cora | 2,708 | 7 | 4,732 | 3,703 |
| Citeseer | 3,327 | 6 | 5,429 | 1,433 |
| Pubmed | 19,717 | 3 | 44,338 | 500 |

For each setting, we conduct 50 random trials, each of which corresponds to a random sampling of the training and attacked nodes and a random initialization of the GCN and the GVAE weights. We compare the accuracy on $V_{attacked}$ before and after copying. "Before copying" refers to the case where we collect the prediction for the attacked nodes from the GCN, which is trained on the attacked graph $\mathcal{G}_{attacked}$. In addition, we report a graph agnostic baseline "Neural Network" on the $V_{attacked}$ set to explore whether it is better to ignore the graph altogether after an attack has been detected.

We employ a Wilcoxon signed-rank test to evaluate the statistical significance of the results obtained. All such tests are performed by comparing with the "Before copying" results. Results marked with an asterisk (*) indicate settings where the test failed to declare a significance at the 5% level.

### A. Ablation Studies

We perform two ablation studies to validate the relevance of the components of the proposed procedure.

*1) Majority voting:* To evaluate the utility of averaging the softmax outputs, we compare with a method when classification is obtained by majority voting. In this method, instead of averaging the softmax outputs at the similar nodes, we make a global decision according to a majority vote among the labels obtained at each nodes. Ties are resolved by random selection.

*2) No copying:* The goal of the second ablation experiment is to ensure that this method is not simply relying on the clustering capability of the chosen embedding technique.
The procedure is the same up to the point where we copy the node. Now, instead of copying the features of the attacked node, we directly take the GCN output of the nodes in $\mathcal{V}_v^{(p)}$ and repeat the two classification procedures: averaging softmax and majority voting.

### B. Results

TABLE II: Average accuracy of the attacked nodes : training with 10 labels per class, $\beta = 50\%$

| Dataset | Before Copying | Copying Average Softmax | Neural Network |
|---|---|---|---|
| Cora | 51.2±8.0 | 56.2±6.2 | 48.9±7.2 |
| Citeseer | 42.2±7.7 | 50.3±7.6 | 39.4±10.0 |
| Pubmed | 51.7±6.7 | 63.3±6.0 | 65.8±6.6 |

TABLE III: Ablation study for majority voting : average accuracy of the attacked nodes for Cora.

| Labels per class | $\beta$ | Before Copying | Copying Majority Voting | Copying Average Softmax |
|---|---|---|---|---|
| 10 | 50% | 51.2±8.0 | 55.7±6.7 | 56.2±6.2 |
| | 75% | 38.8±6.4 | 38.3±8.4 | 39.3±8.0 |
| 20 | 50% | 58.2±6.9 | 58.0±7.5* | 59.1±7.5* |
| | 75% | 32.4±6.4 | 38.3±7.3 | 39.1±7.3 |

TABLE IV: Ablation study for majority voting : average accuracy of the attacked nodes for Citeseer.

| Labels per class | $\beta$ | Before Copying | Copying Majority Voting | Copying Average Softmax |
|---|---|---|---|---|
| 10 | 50% | 42.2±7.7 | 50.0±7.6 | 50.3±7.6 |
| | 75% | 29.0±6.9 | 40.4±6.1 | 40.2±5.9 |
| 20 | 50% | 48.7±5.9 | 52.7±6.3 | 53.1±7.0 |
| | 75% | 31.1±7.0 | 43.8±7.2 | 44.2±6.3 |

TABLE V: Ablation study no copying : Average accuracy of the attacked nodes for Cora

| Labels per class | $\beta$ | No Copying Average Softmax | No Copying Majority Voting | Copying Average Softmax |
|---|---|---|---|---|
| 10 | 50% | 44.2±6.5 | 43.4±6.4 | 56.2±6.2 |
| | 75% | 21.2±6.3 | 21.2±5.8 | 39.3±8.0 |
| 20 | 50% | 45.2±8.6 | 44.7±7.8 | 59.1±7.5* |
| | 75% | 20.1±5.1 | 20.2±5.2 | 39.1±7.3 |

TABLE VI: Ablation study no copying : Average accuracy of the attacked nodes for Citeseer

| Labels per class | $\beta$ | No Copying Average Softmax | No Copying Majority Voting | Copying Average Softmax |
|---|---|---|---|---|
| 10 | 50% | 35.6±6.4 | 35.6±6.4 | 50.3±7.6 |
| | 75% | 21.3±5.6 | 21.1±6.3 | 40.2±5.9 |
| 20 | 50% | 36.8±5.5 | 37.7±5.4 | 53.1±7.0 |
| | 75% | 22.5±5.8 | 22.5±5.6 | 44.2±6.3 |

### C. Discussion

From Tables II, III and IV, we observe that the proposed algorithm offers significant improvement across all datasets from the "Before Copying" baseline at the attacked nodes. In Table II, the relative advantage is apparent as the method is able to improve accuracy by between 5% and 11.6% while outperforming the neural network in most cases. This illustrates the capability of the method of being able to leverage the graph for classification in a situation where labeled data is scarce and not sufficient to train a competitive graph agnostic method.

The ablation studies also confirm that averaging over the softmax performs better than majority voting for almost all experimental settings, but the performance difference is small (Tables III and IV). For the "No Copying" ablation experiment, the results in Tables V and VI show that the proposed copying mechanism offers significant improvement compared to simply using the outputs at the similar nodes. In some cases, "No Copying" is even worse than the performance of "Before Copying", for both softmax averaging and majority voting.

### VI. CONCLUSION

In this paper, we have proposed a recovery algorithm which shows promising results in classifying nodes that have been subjected to a targeted topology attack. The post-attack classification step adds negligible overhead to overall training procedure. We have conducted experiments and ablation studies to highlight the relative importances of different components of the methodology. This work can be further extended by combining the method with an attack detection technique; this will eliminate the assumption that we know the nodes that have been attacked. This scenario is a more realistic setting that we could expect to encounter in practice. In addition, another important research direction is to examine how the depends on the choice of embedding technique and graph-based classifier.

## References

[1] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *Proc. IEEE Symp. Security and Privacy*, San Jose, CA, USA, May 2017, pp. 39–57.

[2] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Proc. Adv. Neural Info. Process. Systems*, Barcelona, Spain, Dec. 2016, pp. 3844–3852.

[3] T. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *Proc. Int. Conf. Learning Representations*, Toulouse, France, Apr. 2017.

[4] W. Hamilton, R. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Proc. Adv. Neural Info. Process. Systems*, Long Beach, CA, USA, Dec. 2017, pp. 1024–1034.

[5] F. Monti, D. Boscaini *et al.*, "Geometric deep learning on graphs and manifolds using mixture model CNNs," in *Proc. IEEE Conf. Comp. Vision and Pattern Recog.*, Honolulu, HI, USA, Jul. 2017, pp. 5425–5434.

[6] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in *Proc. Int. Conf. Machine Learning*, Sydney, Australia, Aug. 2017, pp. 1263–1272.

[7] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *Proc. Int. Conf. Learning Representations*, Vancouver, Canada, Apr. 2018.

[8] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner *et al.*, "Relational inductive biases, deep learning, and graph networks," *arXiv eprint : arXiv:1806.01261*, Oct. 2018.

[9] Y. Zhang, S. Pal, M. Coates, and D. Üstebay, "Bayesian graph convolutional neural networks for semi-supervised classification," in *Proc. AAAI Conf. Artificial Intelligence*, vol. 33, Honolulu, HI, USA, Feb. 2019, pp. 5829–5836.

[10] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec, "Graph convolutional neural networks for web-scale recommender systems," in *Proc. ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, London, UK, Aug. 2018, pp. 974–983.

[11] X. Geng, Y. Li, L. Wang, L. Zhang, Q. Yang, J. Ye, and Y. Liu, "Spatiotemporal multi-graph convolution network for ride-hailing demand forecasting," in *Proc. AAAI Conf. Artificial Intelligence*, vol. 33, Honolulu, HI, USA, Feb. 2019, pp. 3656–3663.

[12] Z. Liu, C. Chen, L. Li, J. Zhou, X. Li, L. Song, and Y. Qi, "Geniepath: Graph neural networks with adaptive receptive paths," in *Proc. AAAI Conf. Artificial Intelligence*, vol. 33, Honolulu, HI, USA, Feb. 2019, pp. 4424–4431.

[13] Z. Deng, Y. Dong, and J. Zhu, "Batch virtual adversarial training for graph convolutional networks," in *Proc. Learning and Reasoning with Graph-Structured Representations Workshop, Intl. Conf. Machine Learning*, Long Beach, CA, Jun. 2019.

[14] K. Sun, P. Koniusz, and J. Wang, "Fisher-bures adversary graph convolutional networks," *arXiv eprint : arXiv 1903.04154*, Jun. 2019.

[15] D. Zügner, A. Akbarnejad, and S. Günnemann, "Adversarial attacks on neural networks for graph data," in *Proc. ACM Int. Conf. Knowl. Disc. Data Mining*, London, UK, Aug. 2018, pp. 2847–2856.

[16] ——, "Adversarial attacks on graph neural networks via meta learning," in *Proc. Int. Conf. Learning Representations*, New Orleans, LA, USA, May 2019.

[17] H. Dai, H. Li, T. Tian, X. Huang, L. Wang, J. Zhu, and L. Song, "Adversarial attack on graph structured data," in *Proc. Int. Conf. Machine Learning*, Stockholm, Sweden, Jul. 2018, pp. 1115–1124.

[18] Y. Zhang, S. Khan, and M. Coates, "Comparing and detecting adversarial attacks for graph deep learning," in *Proc. Representation Learning on Graphs and Manifolds Workshop, Int. Conf. Learning Representations*, New Orleans, LA, USA, May 2019.

[19] T. Kipf and M. Welling, "Variational graph auto-encoders," in *Proc. Bayesian Deep Learning Workshop, Adv. Neural Info. Process. Systems*, Barcelona, Spain, Nov. 2016.

[20] P. Sen, G. Namata *et al.*, "Collective classification in network data," *AI Magazine*, vol. 29, no. 3, p. 93, Sep. 2008.