

# Paper: Mixup for Node and Graph Classification [2]

Jiaxin Liu

Group Reading

July 12, 2021

- 1 Introduction
- 2 Mixup for Graph
- 3 Experiments

- Mix-up example

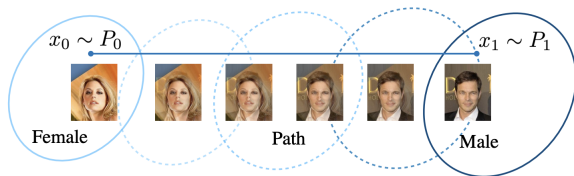


Figure: Example for interpolation.

- Expected Risk:
  - Function  $f : X \mapsto Y$  where  $(x, y) \sim P(X, Y)$ .
  - Loss function  $l : X \times Y \rightarrow \mathbb{R}$
  - $R(f) = \int l(f(x), y) dP(x, y)$

- Expected Risk:

- Function  $f : X \mapsto Y$  where  $(x, y) \sim P(X, Y)$ .
- Loss function  $l : X \times Y \rightarrow \mathbb{R}$
- $R(f) = \int l(f(x), y) dP(x, y)$

- Empirical Risk:

- Training data  $D = \{(x_i, y_i)\}_{i=1}^n$ , where  $(x_i, y_i) \sim P$  for all  $i = 1, \dots, n$ .
- $R_\delta(f) = \int l(f(x), y) dP_\delta(x, y) = \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$

- Pros and Cons? Other methods to approximate  $P$ ? Vicinal Risk Minimization (VRM) [1].

# From VRM to Mixup

- Vicinity distribution  $P_v(\tilde{x}, \tilde{y}) = \frac{1}{n} \sum_{i=1}^n v(\tilde{x}, \tilde{y} | x_i, y_i)$ 
  - $v(\tilde{x}, \tilde{y} | x_i, y_i) = \mathcal{N}(\tilde{x} - x_i, \sigma^2) \delta(\tilde{y} = y_i)$
  - Dataset  $D_v := \{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^m$
  - Empirical vicinal risk:  $R_v(f) = \frac{1}{m} \sum_{i=1}^m l(f(\tilde{x}_i), \tilde{y}_i)$

# From VRM to Mixup

- Vicinity distribution  $P_v(\tilde{x}, \tilde{y}) = \frac{1}{n} \sum_{i=1}^n v(\tilde{x}, \tilde{y} | x_i, y_i)$ 
  - $v(\tilde{x}, \tilde{y} | x_i, y_i) = \mathcal{N}(\tilde{x} - x_i, \sigma^2) \delta(\tilde{y} = y_i)$
  - Dataset  $D_v := \{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^m$
  - Empirical vicinal risk:  $R_v(f) = \frac{1}{m} \sum_{i=1}^m l(f(\tilde{x}_i), \tilde{y}_i)$
- Mixup
  - $\mu(\tilde{x}, \tilde{y} | x_i, y_i) = \frac{1}{n} \sum_j^n \mathbb{E}_\lambda [\delta(\tilde{x} = \lambda \cdot x_i + (1 - \lambda) \cdot x_j, \tilde{y} = \lambda \cdot y_i + (1 - \lambda) \cdot y_j)]$   
where  $\lambda \sim \text{Beta}(\alpha, \alpha)$ , for  $\alpha \in (0, \infty)$ .
  - $\tilde{x} = \lambda x_i + (1 - \lambda) x_j$ ,  $\tilde{y} = \lambda y_i + (1 - \lambda) y_j$  where  $\lambda \in [0, 1]$

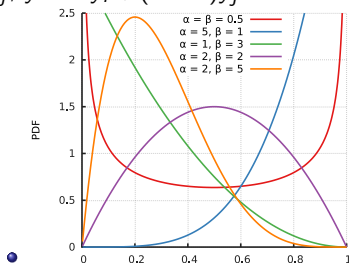


Figure: Beta distributions.

# Mixup Illustration

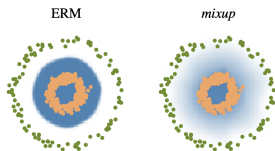


Figure: Toy example. Green: Class 0. Orange: Class 1. Blue shading indicates  $p(y = 1|x)$ . [3]

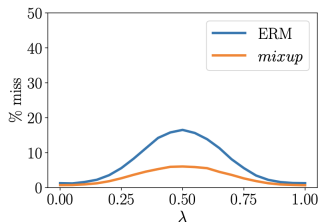


Figure: Prediction error in-between training data. Evaluated at  $x = \lambda x_i + (1 - \lambda)x_j$ , a prediction is counted as a "miss" if it does not belong to  $\{y_i, y_j\}$ . The model trained with mixup has fewer misses. [3]



# Mixup for Graph

- Graph  $G = (V, E)$
- $x_i$ , neighborhood of node  $i$  is  $N(i) = \{j \in V | (i, j) \in E\}$
- GNN:

$$h_i^{(l)} = \text{AGGREGATE}(h_i^{(l-1)}, \{h_j^{(l-1)} | j \in N(i)\}, W^{(l)}).$$

$$h_i^{(0)} = x_i.$$

- Graph classification:  $h_G = \text{READOUT}(\{h_i^{(L)} | i \in V\})$ .

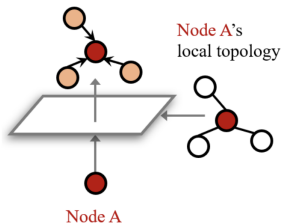


Figure: A GNN layer.

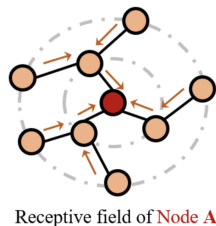


Figure: Receptive field of node A.

# Mixup for Node Classification

- Mixup:  $\tilde{x}_{i,j} = \lambda x_i + (1 - \lambda)x_j$ ,
- Two-branch Mixup for nodes (mix the receptive field subgraphs):

$$\tilde{h}_{ij,i}^{(l)} = \text{AGGREGATE}(\tilde{h}_{ij}^{(l-1)}, \{h_k^{(l-1)} | k \in N(i)\}, W^{(l)})$$

$$\tilde{h}_{ij,j}^{(l)} = \text{AGGREGATE}(\tilde{h}_{ij}^{(l-1)}, \{h_k^{(l-1)} | k \in N(j)\}, W^{(l)})$$

Node mixup

$$\tilde{h}_{ij}^{(l)} = \lambda \tilde{h}_{ij,i}^{(l)} + (1 - \lambda) \tilde{h}_{ij,j}^{(l)}$$

where  $\tilde{h}_{ij}^{(0)} = \tilde{x}_{i,j}$ .

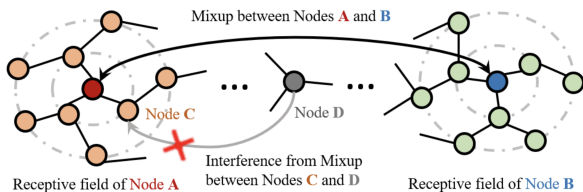


Figure: Two-branch mixup for nodes A and B

# Mixup for Node Classification

How to get  $h_k^{(l)}$  and  $\tilde{h}_{ij}^{(l)}$ ? Two-stage Mixup.

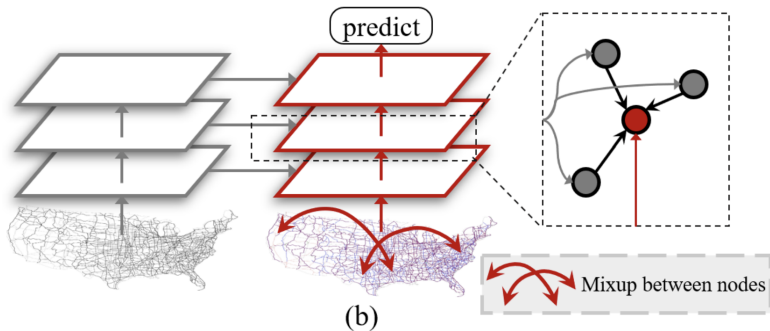


Figure: Two-stage mixup for getting  $h_k^{(l)}$  and  $\tilde{h}_{ij}^{(l)}$ .

# Mixup for Graph Classification

- Only mixup two graphs in the embedding space.

$$\tilde{h}_{G_1, G_2} = \lambda h_{G_1} + (1 - \lambda) h_{G_2}$$

$$\tilde{y}_{G_1, G_2} = \lambda y_{G_1} + (1 - \lambda) y_{G_2}$$

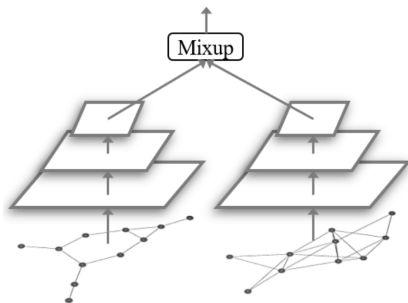


Figure: Mixup for graph classification.

---

**Algorithm 1** Two-Stage Mixup for Node Classification
 

---

**Input:** Graph  $G = (\mathcal{V}, \mathcal{E})$  of a mini-batch, with node attributes  $\{\mathbf{x}_i | i \in \mathcal{V}\}$ , a GNN model with the aggregation function  $\text{AGGREGATE}(\cdot)$ , hyper-parameter  $\alpha$  for the distribution of  $\lambda$ , the ground truth labels  $\{y_i | i \in \mathcal{V}\}$ .

**Output:** The trained parameters of GNN:  $\{\mathbf{W}^{(l)}\}_l$ .

```

1: for  $i \leftarrow 1$  to  $\#\mathcal{V}$  do
2:    $\mathbf{h}_i^{(0)} \leftarrow \mathbf{x}_i$ 
3: end for
4: for  $l \leftarrow 1$  to  $L - 1$  do
5:   for  $i \leftarrow 1$  to  $\#\mathcal{V}$  do
6:      $\mathbf{h}_i^{(l)} \leftarrow \text{AGGREGATE}(\mathbf{h}_i^{(l-1)}, \{\mathbf{h}_j^{(l-1)} | j \in \mathcal{N}(i)\}, \mathbf{W}^{(l)})$ 
7:   end for
8: end for
9: for  $i \leftarrow 1$  to  $\#\mathcal{V}$  do
10:  Sample  $j$  from  $\mathcal{V}$ 
11:   $\lambda \leftarrow \text{Beta}(\alpha, \alpha)$ 
12:   $\tilde{\mathbf{x}}_{ij} \leftarrow \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j$ 
13:   $\tilde{\mathbf{y}}_{ij} \leftarrow \lambda y_i + (1 - \lambda) y_j$ 
14:   $\tilde{\mathbf{h}}_{ij}^{(0)} \leftarrow \tilde{\mathbf{x}}_{ij}$ 
15:  for  $l \leftarrow 1$  to  $L$  do
16:     $\tilde{\mathbf{h}}_{ij,i}^{(l)} \leftarrow \text{AGGREGATE}(\tilde{\mathbf{h}}_{ij}^{(l-1)}, \{\mathbf{h}_k^{(l-1)} | k \in \mathcal{N}(i)\}, \mathbf{W}^{(l)})$   $\longrightarrow$  Branch 1
17:     $\tilde{\mathbf{h}}_{ij,j}^{(l)} \leftarrow \text{AGGREGATE}(\tilde{\mathbf{h}}_{ij}^{(l-1)}, \{\mathbf{h}_k^{(l-1)} | k \in \mathcal{N}(j)\}, \mathbf{W}^{(l)})$   $\longrightarrow$  Branch 2
18:     $\tilde{\mathbf{h}}_{ij}^{(l)} \leftarrow \lambda \tilde{\mathbf{h}}_{ij,i}^{(l)} + (1 - \lambda) \tilde{\mathbf{h}}_{ij,j}^{(l)}$ 
19:  end for
20: end for
21: Calculate classification loss  $\mathcal{L}$  on  $\{\tilde{\mathbf{h}}_{ij}^{(L)}, \tilde{\mathbf{y}}_{ij} | i \in \mathcal{V}\}$ .
22: Back-propagation on  $\{\mathbf{W}^{(l)}\}_l$  for minimizing  $\mathcal{L}$ .
    
```

- Node classification

Method	Citeseer	Cora	Pubmed
GCN [27]	77.1±1.4	88.3±0.8	86.4±1.1
GAT [50]	76.3±0.8	87.6±0.5	85.7±0.7
JKNet [61]	78.1±0.9	89.1±1.2	86.9±1.3
LGCN [18]	77.5±1.1	89.0±1.2	86.5±0.6
GMNN [39]	77.4±1.5	88.7±0.8	86.7±1.0
ResGCN [31]	77.9±0.8	88.1±0.6	87.1±1.2
DropEdge [40] + GCN	78.1±1.1	89.2±0.7	87.3±0.6
DropEdge [40] + JKNet	79.3±0.7	89.9±0.8	87.6±0.9
Mixup + GCN	78.7±0.9	90.0±0.7	87.9±0.8
Mixup + JKNet	<b>80.1±0.8</b>	<b>90.4±0.9</b>	<b>88.3±0.6</b>

Figure: Test Accuracy of transductive node classification.

- Varied training ratios  $r$ .

Method	Citeseer			Cora			Pubmed		
	$r = 30\%$	$r = 40\%$	$r = 50\%$	$r = 30\%$	$r = 40\%$	$r = 50\%$	$r = 30\%$	$r = 40\%$	$r = 50\%$
GCN [27]	$74.7 \pm 2.5$	$75.2 \pm 1.8$	$76.3 \pm 1.6$	$86.3 \pm 1.9$	$86.8 \pm 1.4$	$87.5 \pm 1.0$	$85.1 \pm 2.3$	$85.4 \pm 1.4$	$85.8 \pm 1.2$
Mixup + GCN	$76.9 \pm 2.1$	$77.1 \pm 1.5$	$78.1 \pm 1.3$	$88.5 \pm 1.4$	$88.9 \pm 1.0$	$89.4 \pm 0.9$	$87.0 \pm 1.6$	$87.2 \pm 1.1$	$87.5 \pm 1.0$
JKNet [61]	$75.6 \pm 1.9$	$76.0 \pm 1.4$	$77.1 \pm 1.1$	$86.7 \pm 2.1$	$87.4 \pm 1.5$	$88.2 \pm 1.3$	$85.3 \pm 2.2$	$85.9 \pm 1.6$	$86.4 \pm 1.4$
Mixup + JKNet	<b><math>78.0 \pm 1.7</math></b>	<b><math>78.3 \pm 1.2</math></b>	<b><math>79.2 \pm 1.0</math></b>	<b><math>88.6 \pm 2.0</math></b>	<b><math>89.1 \pm 1.5</math></b>	<b><math>89.7 \pm 1.2</math></b>	<b><math>87.2 \pm 1.9</math></b>	<b><math>87.5 \pm 1.3</math></b>	<b><math>87.9 \pm 0.9</math></b>

Figure: Test Accuracy of node classification over different training set ratios.

# Experiments

- t-SNE plot for the final-layer representations.
- Loss on the test data during training.

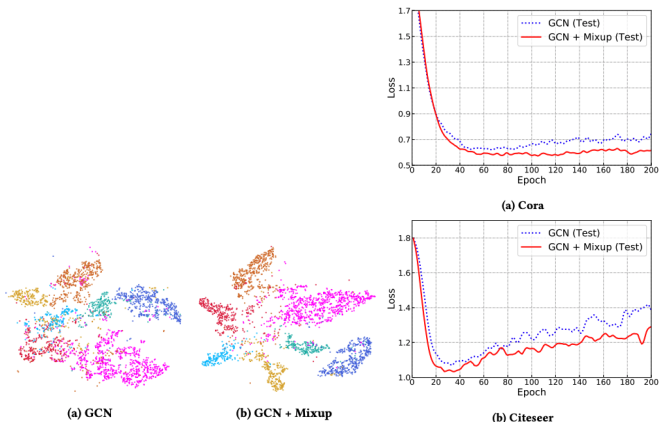


Figure: The learned representations of the nodes in the Cora dataset.

Figure: The training curves of GCN with and without Mixup.



# Experiments

- Two-stage framework.

Method	two stages	Pubmed	$\Delta$	Yelp	$\Delta$
GCN [27]	-	86.4 $\pm$ 1.1	0	65.3 $\pm$ 0.3	0
Mixup + GCN	w/o	85.8 $\pm$ 1.3	-0.6	64.2 $\pm$ 0.6	-1.1
	w/	87.9 $\pm$ 0.8	+1.5	66.3 $\pm$ 0.4	+1.0

Figure: The node classification results with and without two-stage framework.

- Selection for  $\alpha$  in Beta( $\alpha, \alpha$ ).

Method	$\alpha$	Citeseer	Cora	Flickr
GCN [27]	-	77.1 $\pm$ 1.4	88.3 $\pm$ 0.8	51.1 $\pm$ 0.2
Mixup + GCN	0.2	78.1 $\pm$ 0.9	89.2 $\pm$ 0.8	52.0 $\pm$ 0.3
	0.5	78.4 $\pm$ 0.8	89.5 $\pm$ 0.7	52.1 $\pm$ 0.3
	1	<b>78.7<math>\pm</math>0.9</b>	<b>90.0<math>\pm</math>0.7</b>	52.4 $\pm$ 0.4
	2	78.6 $\pm$ 1.0	89.8 $\pm$ 0.8	<b>52.8<math>\pm</math>0.5</b>
	5	78.4 $\pm$ 1.2	89.4 $\pm$ 1.1	52.7 $\pm$ 0.4

Figure: Node classification results with different



Olivier Chapelle et al. “Vicinal risk minimization”. In: *Advances in neural information processing systems* (2001), pp. 416–422.



Yiwei Wang et al. “Mixup for Node and Graph Classification”. In: *Proceedings of the Web Conference 2021*. 2021, pp. 3663–3674.



Hongyi Zhang et al. “mixup: Beyond empirical risk minimization”. In: *arXiv preprint arXiv:1710.09412* (2017).