

GNN Fundamentals - 3

In previous part - we have seen how can we transform the output of GNN embedding to various tasks like node classification, edge classification. Here we will see Graph level tasks.

Graph level task heard?

$$\hat{y}_G = \text{Headgraph} \left(h_v \in \mathbb{R}^d, \forall v \in G \right)$$

We need to design Headgraph which will aggregate each node embedding. The aggregated node embedding can be used for downstream task.

Suppose graph has K-class problem where we need to do K-way prediction.

$$\hat{y}_K = \text{Transform}_{\text{graph}} \left(\hat{y}_G \in \mathbb{R}^d \right)$$

$$\hat{y}_K \in \mathbb{R}^K$$

There will be many ways in which we can Aggregate

$$h_v^L$$

① Global pooling.

② Global mean pooling: Mean($h_v \in \mathbb{R}^d, \forall v \in G$)

↳ coordinate wise take mean.

③ Global max pooling: Max($h_v \in \mathbb{R}^d, \forall v \in G$)

↳ coordinate wise take max value.

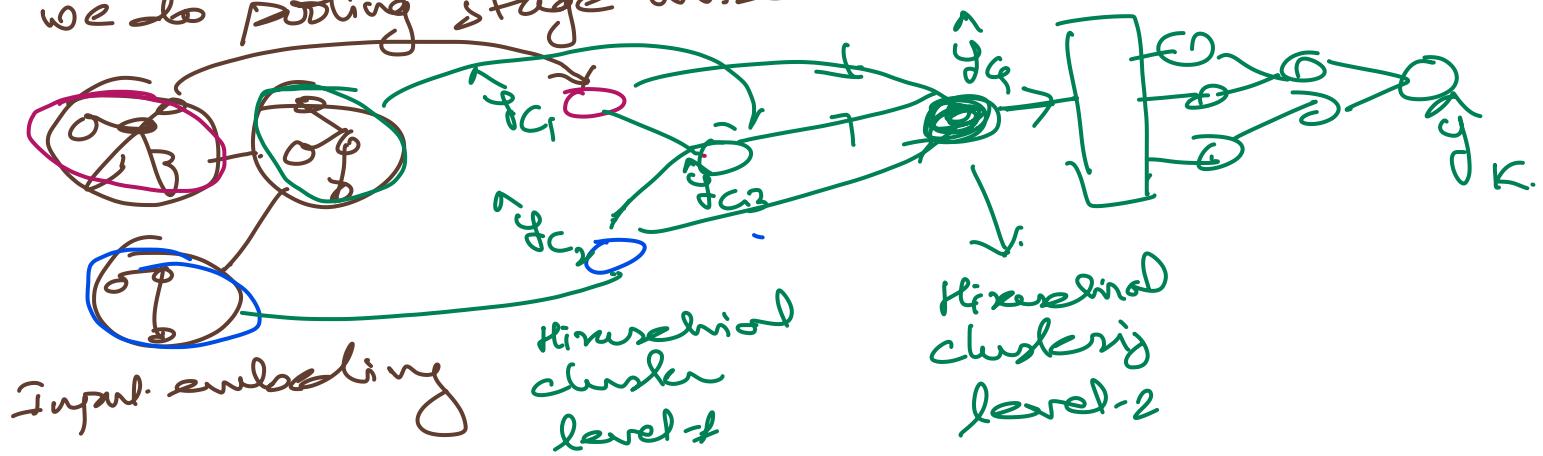
④ Global sum pooling: sum($h_v \in \mathbb{R}^d, \forall v \in G$)

These option works well with small graph but - not for large graph as when we do Global pooling then many details at graph level lost. : 1-dimensionality.

e.g.: $G_1 = \hat{y}_{G_1} = \{-1, -2, 0, 1, 2\}$, $G_2 \Rightarrow \hat{y}_{G_2} = \{-10, -20, 0, 10, 20\}$

If we do sum($\hat{y}_{G_1}\right) = 0 + \sum(\hat{y}_{G_2}) = 0$ but both graph node embedding is different - so graph

will not be same.
 For large graphs Global pooling is not effective so we do hierarchical pooling. For this we cluster the similar node embedding and for each cluster we do pooling stage wise.



Two GNN used

GNN-A: Compute node embedding

GNN-B: cluster similar node embeddings. This is also

hierarchical aggregation of node embedding.

Finally single node embedding receives all-to-all
hierarchical and used for downstream tasks.

Now we will talk about - prediction and labels.

When we take ground truth from external source
then that is supervised learning e.g taking class label
value from external source.

when supervision is missing

graph like graph signal for clustering that is called
unsupervised or self-supervised e.g predicting node
clustering coefficient or check if two graphs is
isomorphic. For this we do not take any external
reference for label, so this is unsupervised learning.

Supervised label: This comes from specific use case like in citation network node label y_v will be which subject the node belongs to.

In case of transaction network edge label could be fraud or not. In case of molecular graph if molecule is toxic or toxicity of molecule.

If we can transform the actual task to node, link or graph level prediction task then we can reuse lot of research model in this area.

Similarly for unsupervised task like predict pageRank, clustering coefficient, predict missing link or what kind of molecule is this? for all this we do not need ground truth but we use graph signal itself to predict these.

Now we saw prediction and labels. Now

Till now we saw prediction and labels. Now we will evaluate discrepancy between prediction and label. For this we define loss function and optimization of loss function will do back propagation to optimize all learnable parameters. Evaluation matrix will be used to check performance of model.

Setting up for GNN training:

Prediction	Ground truth	N-data point - can be for any thing so we say \hat{y}_i & y_i as datapoint applicable to all-
For node: \hat{y}_v	y_v	
For edge: \hat{y}_{uv}	y_{uv}	
For graph: \hat{y}_G	y_G	

class label y_i can be for classification or regression.
In classification y_i will be categorical value (e.g. in regression it will be continuous value)
eg whether molecule is toxic or not [Binary classification]
or what is toxicity of molecule [Regression].
Based on type of problem we define loss function.

Classification loss: Most popular loss function is cross entropy to estimate loss when we do classification.

$$CE(\hat{y}_i, y_i) = - \sum_{j=1}^k y_j \log(\hat{y}_j)$$

i is i^{th} data point, j = class of k -class model.

so $CE(\hat{y}_i, y_i)$ represent cross entropy loss for i^{th} data point
For N data point need to sum $\sum_{i=1}^N CE(\hat{y}_i, y_i)$.

For regression kind of task: Suppose we have k -way regression then we can use mean-square error loss (MSE) for each data point i^{th} :

$$MSE(\hat{y}_j^{(i)}, y_j^{(i)}) = \sum_{j=1}^k (\hat{y}_j^{(i)} - y_j^{(i)})^2$$

For all data point-
 $Loss = \sum_{i=1}^N MSE(\hat{y}_j^{(i)}, y_j^{(i)}) \Rightarrow \min_{\theta} (Loss) \rightarrow$ objective function

when loss function optimizes all parameter then we need metrics for evaluation of performance of the model. We use standard evaluation matrix for GNN like confusion matrix & ROC (Recevier operating curve)

In case of binary classification:

$$\text{Accuracy} = \frac{TP+TN}{\text{Total data set}}$$

$$(P) \text{Precision} = \frac{TP}{TP+FP} = \frac{\text{Actual true out of model predicted true}}{\text{Total Model predicted true class}}$$

$$(R) \text{Recall} = \frac{\text{Actual true out of model predicted true}}{\text{Total true class.}} = \frac{TP}{TP+FN}$$

$$F1\text{-score}(F) = \frac{1}{F} = \frac{1}{P} + \frac{1}{R}$$

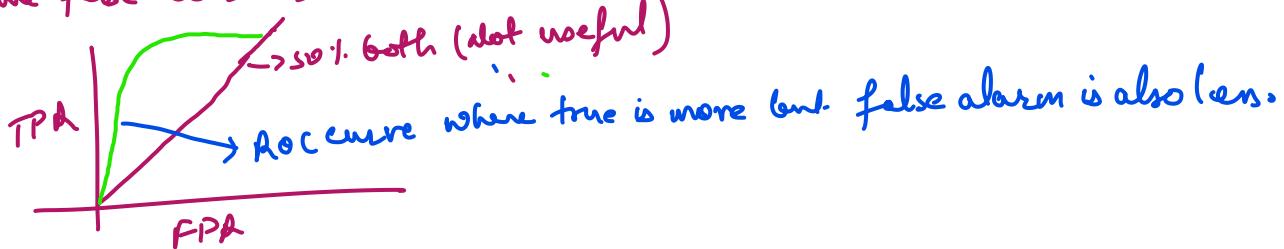
AUC Curve: This is trade-off between true alarm & false alarm.

TP is part of true alarm & FP is part of false alarm.

$$TPA(\text{True positive rate}) = \frac{TP}{TP+FN} = \frac{\text{Actual true out of model predicted true}}{\text{Total actual true class.}}$$

$$FPA(\text{False positive rate}) = \frac{FP}{FP+TN} = \frac{\text{Model predicted positive which is false}}{\text{Total negative class.}}$$

So for any model it is important to have high true alarm but at the same time false alarm should be minimum. This is represented as curve (ROC)



Data split for training, evaluation and test.

