

Learning how to use GNN

To get started, you need to install pytorch geometric:

<https://pytorch-geometric.readthedocs.io/en/latest/install/installation.html>.

You can then have a look at the tutorial:

https://pytorch-geometric.readthedocs.io/en/latest/get_started/introduction.html.

Check in particular some sections:

- Data Handling of Graphs
- Data Transforms
- Learning Methods on Graphs

1. Getting started: data preparation for node classification

- Load the toy dataset `ToyFriendship.graphml` from the class website, using `networkx`. Plot the network to have a quick view, check the attributes (`G.nodes(data=True)`). We consider that we know the preferences of the students among sports/music/science, and the club they belong to.
- Convert from `networkx` to pytorch geometric using `torch_geometric.utils.from_networkx` function. Be careful, due to some bug, you first need to do `G.graph={}` on your `networkx` graph. In the function, use `group_node_attrs=["like_sports","like_music","like_science"]` to load only those attributes as `x`.
- Check what is inside this object. You should find the edges `edge_index`, the node features `x`, etc.
- Encode the class (`club`) using sklearn `LabelEncoder`, e.g.,

```
encoder = LabelEncoder()
integer_labels = encoder.fit_transform(data.club)
target_tensor = torch.tensor(integer_labels, dtype=torch.long)
data.y = target_tensor
data.num_classes = len(set(data.club))
```

- Let's consider that for some of the students, we don't know their preferences, but we want to train a model to guess the club they belong to. For instance, we can imagine new students to whom we want to recommend a club. So we want to guess the club class from the `like` attributes, but for students for which we don't have the like attribute. Without graph, this is not possible.
- We need to hide the `like` information for some of the nodes. You can do it with a mask, with something like:

```

num_nodes = data.num_nodes
train_ratio = 0.80 # 80% of nodes for training

# Randomly creating a mask
mask = torch.rand(num_nodes) < train_ratio
data.train_mask = mask
data.test_mask = ~data.train_mask

# remove the attributes for the nodes that are not in the training set
temp = torch.zeros((num_nodes, 3), dtype=torch.float)
temp[data.train_mask] = data.x[data.train_mask]
data.x = temp

```

2. Predict using a GCN

- (a) Build your first GCN, with a single layer. It should solve a classification problem, with 3 classes.
- (b) Your evaluation should be only on the test set, i.e., something like:

```

pred = model(data).argmax(dim=1)
correct = (pred[data.test_mask] == data.y[data.test_mask]).sum()
acc = int(correct) / int(data.test_mask.sum())
print(f'Accuracy: {acc:.4f}')

```

- (c) Check the `add_self_loops` attributes of the conv layer, and think of its meaning.
- (d) Print the targets and the predictions for all the nodes
- (e) Plot a confusion matrix, e.g.,

```

from sklearn.metrics import confusion_matrix
import seaborn as sns
confmat = confusion_matrix(data.y[data.test_mask], pred[data.test_mask])
sns.heatmap(confmat, annot=True, fmt='g')

```

- (f) Print the weights of the GCN layer and interpret them (if you have a good accuracy...)

3. Predicting edges using a VGAE

- (a) This time, we want to predict edges. Start back from the original network, without hidden information. Build a train_test split using the `train_test_split_edges` function from `torch_geometric.utils`. You don't need a validation set.
- (b) Build your `Encoder`. It should be something like:

```

class Encoder(torch.nn.Module):
    def __init__(self, in_channels, out_channels):
        super().__init__()
        self.conv1 = GCNConv(in_channels, 2*out_channels)
        self.conv_mu = GCNConv(2*out_channels, out_channels)
        self.conv_logstd = GCNConv(2*out_channels, out_channels)

    def forward(self, x, edge_index):
        x = self.conv1(x, edge_index).relu()
        return self.conv_mu(x, edge_index), self.conv_logstd(x, edge_index)

```

(The VGAE needs to outputs, mu and logstd, which corresponds for each encoded node to a centroid and a standard deviation. Each dot is encoded as a probability, to avoid "overfitting" each point in a particular position. Check more details in the class on the VGAE to know more.

- (c) Initialize your model using `VGAE` from `torch_geometric.nn`.
- (d) Evaluate your model. Be careful to use training data for training and testing data for testing :) Something like:

```

z = model.encode(data.x, data.train_pos_edge_index)
return model.test(z, data.test_pos_edge_index, data.test_neg_edge_index)

```

- (e) Evaluate the result of your test, i.e., with AUC and AP (Average Precision)
- (f) Check that you are able to reconstruct the original graph, by applying the dot product between node vectors. You can do it with something like:

```

z = model.encode(data.x, data.train_pos_edge_index)
Ahat = torch.sigmoid(z @ z.T)

```

1 Going Further

- (a) If you want to do more, you can use the airport dataset with population provided in the class' page. Can you predict edges on that dataset? What about predicting the population? The country?