

# Data Science Capstone

## Graph-Based Neural Networks (GNN)

### Recommendation Systems



THE GEORGE  
WASHINGTON  
UNIVERSITY

WASHINGTON, DC

# A Pipeline of Bipartite Graph Neural Network Operator for Recommendation Systems Using PyTorch Geometric



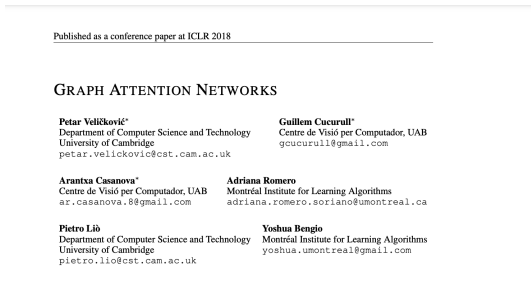
- **Heterogeneous Graph Learning** refers to any graph-based neural network learning on a graph with more than two vertex or edge types.
- Standard Message Passing GNNs (MP-GNNs) can not trivially be applied to heterogeneous graph data, as vertex and edge features from different types can not be processed by the same functions due to differences in feature type.
- To automatically convert a homogeneous GNN model to a heterogeneous GNN model by making use of:

```
1 | from torch_geometric.nn import to_hetero
```

# Case Study: Heterogeneous Graph Learning

# Graph Attention Network (GAT)

- The "Graph Attention Networks" research paper was first published in October 30th, 2017.



- The main idea is to compute the hidden representations of each vertex in the graph, by attending over its neighbors, following a self-attention strategy.

# (1) GAT Layer: Linear Layer

- We begin by applying a shared linear transformation, parameterized by a weight matrix  $\mathbf{W}$  on all of the vertex features  $\mathbf{h} = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_N\}$ :

$$\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_j$$

# (1) Case Study: Linear Layer

- For recommendations systems, we apply two distinct shared linear transformations, parameterized by weight matrices  $\mathbf{W}_{\text{user}}$  and  $\mathbf{W}_{\text{item}}$  on the sampled sets of bipartite vertex features  $\mathbf{h}_{\text{user}} = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_N\}$  and  $\mathbf{h}_{\text{item}} = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_M\}$ :

$$\mathbf{W}_{\text{user}} \vec{h}_{\text{user}_i}, \mathbf{W}_{\text{user}} \vec{h}_{\text{user}_j}$$

$$\mathbf{W}_{\text{item}} \vec{h}_{\text{item}_i}, \mathbf{W}_{\text{item}} \vec{h}_{\text{item}_j}$$



# (1) Linear Layer: Code

- Initializing  $W_{\text{user}}$  and  $W_{\text{item}}$ :

```
1 | self.lin_src = Linear(in_channels[0], heads * out_channels, False,  
2 |                       weight_initializer='glorot')  
3 | self.lin_dst = Linear(in_channels[1], heads * out_channels, False,  
4 |                       weight_initializer='glorot')
```

- Applying linear transformation and reshaping into a tensor of size (Batch, Heads, Output Channels):

```
1 | # If the module is initialized as bipartite, transform source  
2 | # and destination node features separately:  
3 | assert self.lin_src is not None and self.lin_dst is not None  
4 | x_src = self.lin_src(x).view(-1, H, C)  
5 | x_dst = self.lin_dst(x).view(-1, H, C)
```

- Note: 'Src' and 'Dst' refer to Source and Destination vertex features; in our bipartite graph these are user and item vertex features.

## (2) GAT Layer: Self-Attention

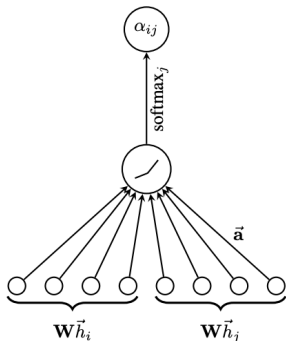
- Then, we perform self-attention on the vertices.

$$\alpha_{i,j} = \frac{\exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}^\top [\mathbf{W}\vec{h}_i \parallel \mathbf{W}\vec{h}_j]\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}^\top [\mathbf{W}\vec{h}_i \parallel \mathbf{W}\vec{h}_k]\right)\right)}$$

- The attention coefficients  $\alpha_{i,j}$  indicate the importance of vertex  $j$ 's features to vertex  $i$ .

## (2) GAT Layer: Self-Attention

- The attention mechanism  $\vec{a}$  is a learnable parameter that is multiplied to the concatenation of the output of the previous linear layer.
- The output is then passed through a LeakyReLU activation with negative slope 0.2, and a softmax normalization across all neighbors.



## (2) Case Study: Self-Attention

- For recommendation systems, we have two different attention coefficients calculated in the same manner:

$$\alpha_{\text{user}_{i,j}} = \frac{\exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}_{\text{user}}^\top [\mathbf{W}_{\text{user}} \vec{h}_{\text{user}_i} \parallel \mathbf{W}_{\text{user}} \vec{h}_{\text{user}_j}]\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}_{\text{user}}^\top [\mathbf{W}_{\text{user}} \vec{h}_{\text{user}_i} \parallel \mathbf{W}_{\text{user}} \vec{h}_{\text{user}_j}]\right)\right)}$$

$$\alpha_{\text{item}_{i,j}} = \frac{\exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}_{\text{item}}^\top [\mathbf{W}_{\text{item}} \vec{h}_{\text{item}_i} \parallel \mathbf{W}_{\text{item}} \vec{h}_{\text{item}_j}]\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}_{\text{item}}^\top [\mathbf{W}_{\text{item}} \vec{h}_{\text{item}_i} \parallel \mathbf{W}_{\text{item}} \vec{h}_{\text{item}_j}]\right)\right)}$$

## (2) Self-Attention: Code

- Initializing  $\vec{a}_{\text{user}}$  and  $\vec{a}_{\text{item}}$ :

```
1 | # The learnable parameters to compute attention coefficients:
2 | self.att_src = Parameter(torch.empty(1, heads, out_channels))
3 | self.att_dst = Parameter(torch.empty(1, heads, out_channels))
```

- These parameters are then passed through glorot intializer:

```
1 | glorot(self.att_src)
2 | glorot(self.att_dst)
```

- Applying self-attention mechanism on previous output:

```
1 | # Next, we compute node-level attention coefficients, both for source
2 | # and target nodes (if present):
3 | alpha_src = (x_src * self.att_src).sum(dim=-1)
4 | alpha_dst = None if x_dst is None else (x_dst * self.att_dst).sum(-1)
5 | alpha = (alpha_src, alpha_dst)
```

## (2) Self-Attention: Code

- Sampled self-attention coefficients are collected according to edge index:

$\alpha_{\text{user}_{i,j}}$

```
1 | # Collect user-defined arguments:
2 |
3 | # (1) - Collect 'alpha_j':
4 | if isinstance(alpha, (tuple, list)):
5 |     assert len(alpha) == 2
6 |     _alpha_0, _alpha_1 = alpha[0], alpha[1]
7 |     if isinstance(_alpha_0, Tensor):
8 |         self._set_size(size, 0, _alpha_0)
9 |         alpha_j = self._index_select(_alpha_0, edge_index_j)
10 |     else:
11 |         alpha_j = None
12 |     if isinstance(_alpha_1, Tensor):
13 |         self._set_size(size, 1, _alpha_1)
14 |     elif isinstance(alpha, Tensor):
15 |         self._set_size(size, j, alpha)
16 |         alpha_j = self._index_select(alpha, edge_index_j)
```

$\alpha_{\text{item}_{i,j}}$

```
1 | # (2) - Collect 'alpha_i':
2 | if isinstance(alpha, (tuple, list)):
3 |     assert len(alpha) == 2
4 |     _alpha_0, _alpha_1 = alpha[0], alpha[1]
5 |     if isinstance(_alpha_0, Tensor):
6 |         self._set_size(size, 0, _alpha_0)
7 |     if isinstance(_alpha_1, Tensor):
8 |         self._set_size(size, 1, _alpha_1)
9 |         alpha_i = self._index_select(_alpha_1, edge_index_i)
10 |     else:
11 |         alpha_i = None
12 | elif isinstance(alpha, Tensor):
13 |     self._set_size(size, i, alpha)
14 |     alpha_i = self._index_select(alpha, edge_index_i)
```

## (2) Self-Attention: Code

- These sampled attention coefficients are summed  $\alpha_{user_{i,j}} + \alpha_{item_{i,j}}$  to emulate concatenation operator  $\parallel$

```
1 | # Given edge-level attention coefficients for source and target nodes,  
2 | # we simply need to sum them up to "emulate" concatenation:  
3 | alpha = alpha_j if alpha_i is None else alpha_j + alpha_i
```

- LeakyReLU activation with negative slope 0.2, and a softmax normalization across all neighbors:

```
1 | alpha = F.leaky_relu(alpha, self.negative_slope)  
2 | alpha = softmax(alpha, index, ptr, dim_size)
```

### (3) GAT Layer: Aggregation

- Once obtained, the normalized attention coefficients are used to compute a linear combination of the features corresponding to them, to serve as the final output features for every vertex:

$$\vec{h}_i' = \sum_{j \in \mathcal{N}_i} \alpha_{i,j} \mathbf{W} \vec{h}_j$$



### (3) Aggregation: Case Study

- For recommendation systems, only user final outputs vertex features  $\vec{h}'_{\text{user}_i}$  are propagated:

$$\vec{h}'_{\text{user}_i} = \sum_{j \in \mathcal{N}_i} (\alpha_{\text{user}_{i,j}} + \alpha_{\text{item}_{i,j}}) \mathbf{W}_{\text{user}} \vec{h}_{\text{user}_j}$$

### (3) Aggregation: Code

- User neighborhood vertices are collected by sampled edge index:

$$\mathbf{W}_{\text{user}} \vec{h}_{\text{user}_j}$$

```
1 | # Collect user-defined arguments:
2 | # (1) - Collect 'x_j':
3 | if isinstance(x, (tuple, list)):
4 |     assert len(x) == 2
5 |     _x_0, _x_1 = x[0], x[1]
6 |     if isinstance(_x_0, Tensor):
7 |         self._set_size(size, 0, _x_0)
8 |         x_j = self._index_select(_x_0, edge_index_j)
9 |     else:
10 |         x_j = None
11 |     if isinstance(_x_1, Tensor):
12 |         self._set_size(size, 1, _x_1)
13 | elif isinstance(x, Tensor):
14 |     self._set_size(size, j, x)
15 |     x_j = self._index_select(x, edge_index_j)
16 | else:
17 |     x_j = None
```

### (3) Aggregation: Code

- Self-attention coefficient  $\alpha_{user_{i,j}} + \alpha_{item_{i,j}}$  is multiplied to the linear layer output  $\mathbf{W}_{user} \vec{h}_{user_j}$

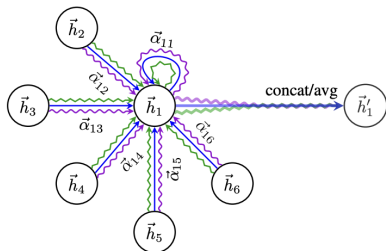
```
1 | def message(self, x_j: Tensor, alpha: Tensor) -> Tensor:
2 |     return alpha.unsqueeze(-1) * x_j
```

- Lastly, the message is propagated to neighborhood vertices  $\sum_{j \in \mathcal{N}_i}$

```
1 | base_cls = aggr.Aggregation
2 | aggrs = [
3 |     aggr for aggr in vars(aggr).values()
4 |     if isinstance(aggr, type) and issubclass(aggr, base_cls)
5 | ]
6 | aggr_dict = {
7 |     'add': aggr.SumAggregation,
8 | }
9 | return resolver(aggrs, aggr_dict, query, base_cls, None, *args, **kwargs)
```

## (Optional 4) GAT Layer: Multi-Head Attention

- Multi-head attention is an optional addition to stabilize the learning process of self-attention:



$$\vec{h}'_i = \frac{1}{k} \sum_{k=1}^k \sum_{j \in \mathcal{N}_i} \alpha_{i,j} \mathbf{W} \vec{h}_j$$

- We simply repeat the process  $k$  times and average the messages.

## (Optional 4) Multi-Head Attention: Case Study

- For recommendation systems, only user final outputs vertex features  $\vec{h}'_{user_i}$  are propagated and averaged:

$$\vec{h}'_{user_i} = \frac{1}{k} \sum_{k=1}^k \sum_{j \in \mathcal{N}_i} (\alpha_{user_i,j} + \alpha_{item_i,j}) \mathbf{W}_{user} \vec{h}_{user_j}$$

## (Optional 4) Multi-Head Attention: Code

- To establish multi-head attention, one must set the 'concat' to False and specify the amount of heads:

```
1 |         if self.concat:
2 |             out = out.view(-1, self.heads * self.out_channels)
3 |         else:
4 |             out = out.mean(dim=1)
```

```
1 | class Encoder(torch.nn.Module):
2 |
3 |     def __init__(self, hidden_channels, heads):
4 |         super().__init__()
5 |         self.conv1 = GATConv(in_channels = (-1, -1), out_channels = hidden_channels, add_self_loops=False, heads = heads, concat=False)
6 |         self.lin1 = Linear(-1, hidden_channels)
7 |         self.conv2 = GATConv(in_channels = (-1, -1), out_channels = hidden_channels, add_self_loops=False, heads = heads, concat=False)
8 |         self.lin2 = Linear(-1, hidden_channels)
9 |         self.conv3 = GATConv(in_channels = (-1, -1), out_channels = hidden_channels, add_self_loops=False, heads = heads, concat=False)
10 |         self.lin3 = Linear(-1, hidden_channels)
11 |
12 |     def forward(self, batch_embedding_dict, batch_edge_index_dict):
13 |
14 |         batch_embedding_dict = self.conv1(batch_embedding_dict, batch_edge_index_dict) + self.lin1(batch_embedding_dict)
15 |         batch_embedding_dict = batch_embedding_dict.relu()
16 |         batch_embedding_dict = self.conv2(batch_embedding_dict, batch_edge_index_dict) + self.lin2(batch_embedding_dict)
17 |         batch_embedding_dict = batch_embedding_dict.relu()
18 |         batch_embedding_dict = self.conv3(batch_embedding_dict, batch_edge_index_dict) + self.lin3(batch_embedding_dict)
19 |         batch_embedding_dict = batch_embedding_dict.relu()
20 |
21 |         return batch_embedding_dict
```

- Once we have our output vertex features for both user and item, we can decode them by performing a Hadamard multiplication and averaging across the  $N$  column dimension to get output  $\vec{d}$ :

$$\vec{d} = \frac{1}{N} \sum_{N=1}^N \mathbf{h}'_{\text{user}} \circ \mathbf{h}'_{\text{item}}$$



# Inner Product Decoder: Code

```
1 | class Decoder(torch.nn.Module):  
2 |  
3 |     def forward(self, batch_embedding_dict, edge_index):  
4 |  
5 |         user_embedding = batch_embedding_dict['user'][edge_index[0]]  
6 |         item_embedding = batch_embedding_dict['item'][edge_index[1]]  
7 |  
8 |         return (user_embedding * item_embedding).sum(dim=-1)
```

- Finally, now that we have our output vector  $\vec{d}$  we can calculate a mean squared error (MSE) with our target edge label vector  $\vec{t}$ :

$$e = \frac{1}{n} \sum_{i=1}^n (\vec{t} - \vec{d})^2$$

- In this manner, we can approximate edge label values to predict the existence of edges between users and items, or vice versa.

# Loss Function: Code

```
1 | out = model(batch.x_dict,  
2 |             batch.edge_index_dict,  
3 |             batch['user', 'item'].edge_label_index)  
4 | y = batch['user', 'item'].edge_label  
5 |  
6 | loss = torch.nn.functional.mse_loss(out, y.float())
```



# Maximum Inner Product Search (MIPS)

- Given a collection of “database” vectors,  
 $D = \{\vec{x}_0, \vec{x}_1, \dots, \vec{x}_d\} \in S \subset R^d$  and a query  $\vec{q} \in R^d$ ,  
**Maximum Inner Product Search** attempts to find a data vector  $\vec{x}_*$  maximizing the inner product with the query:

$$\vec{x}_* = \operatorname{argmax}_{\vec{x} \in S} \vec{q}^T \vec{x}$$

- There are many variations of MIPS, one of them being **Nearest Neighbor MIPS** attempting to find a data vector  $\vec{x}_*$  by minimizing over a distance function  $\rho$ :

$$\vec{x}_* = \operatorname{argmin}_{\vec{x} \in S} \rho(\vec{q}^T \vec{x})$$

Intuition: A neighbor of my neighbor is likely to be my neighbor.

# K-Nearest Neighbor MIPS: Pseudocode

---

## Algorithm KNN-MIPS

---

**Input:** a KNN graph  $G = (D, E)$ ; a query point  $\vec{q}$ ; the number of required nearest neighbors  $K$ ; the number of random restarts  $R$ ; distance function  $\rho$ ; the number of greedy steps  $T$ ; the number of expansions  $E$ .

$S \leftarrow \{\}$  {Set of MIPS candidates}  
 $U \leftarrow \{\}$  {Set of to be expanded candidates}  
**for**  $r = 1, \dots, R$  **do**  
     $\vec{x}_0$  a point drawn randomly from a uniform distribution over  $D$   
    **for**  $t = 1, \dots, T$  **do**  
         $\vec{x}_t = \operatorname{argmin}_{\vec{x} \in N(\vec{x}_{t-1}, E, G)} \rho(\vec{q}^T \vec{x})$   
         $S \leftarrow N(\vec{x}_{t-1}, E, G)$  { $N$  is a function that returns neighbors}  
         $U \leftarrow \{\rho(\vec{q}^T \vec{x}) : \vec{x} \in N(\vec{x}_{t-1}, E, G)\}$ .  
    **end for**  
**end for**

**Output:** Sort  $U$ , pick the first  $K$  elements, and return the corresponding indices in  $S$ .

---

# K-Nearest Neighbor MIPS: Visualization

- Suppose greedy steps  $T = 2$ , neighbors  $K = 1$  and expansions  $E = 3$



# K-Nearest Neighbors MIPS: Code

- At inference, we first calculate the embedding, or output features, for item vertices  $\mathbf{h}'_{\text{item}}$

```
1 | positive_item_embeddings = []
2 | for batch in item_loader:
3 |     batch = batch.to(device)
4 |     batch_embedding_item = model.encoder(batch.x_dict,
5 |                                         batch.edge_index_dict)['item']
6 |     batch_embedding_item_positive = batch_embedding_item[:batch['item'].batch_size]
7 |
8 |     positive_item_embeddings.append(batch_embedding_item_positive)
9 | horizontal_stack_positive_item_embeddings = torch.cat(positive_item_embeddings, dim = 0)
10 | del positive_item_embeddings
11 |
12 | max_inner_product_search = MIPSKNNIndex(horizontal_stack_positive_item_embeddings)
```

- Then, we pass these item embedding vertices  $\mathbf{h}'_{\text{item}}$  as our dataset vectors  $D = \{\vec{x}_0, \vec{x}_1, \dots, \vec{x}_d\}$
- Note: positive embedding are passed as dataset vectors if we are utilizing negative sampling.

# K-Nearest Neighbors MIPS: Code

- After passing dataset vectors, we calculate the embedding, or output features, for user vertices  $\mathbf{h}'_{\text{user}}$

```
1 | for batch in user_loader:
2 |     batch = batch.to(device)
3 |     batch_embedding_user = model.encoder(batch.x_dict,
4 |                                         batch.edge_index_dict)['user']
5 |     batch_embedding_user_positive = batch_embedding_user[:batch['user'].batch_size]
6 |
7 |     edge_index = test_edge_label_index.sparse_narrow(
8 |         dim=0,
9 |         start=number_processed,
10 |         length=batch_embedding_user_positive.size(0))
11 |
12 |     exclude_links = test_exclude_links.sparse_narrow(
13 |         dim=0,
14 |         start=number_processed,
15 |         length=batch_embedding_user_positive.size(0))
16 |
17 |     number_processed += batch_embedding_item_positive.size(0)
18 |
19 |     _, pred_index_mat = max_inner_product_search.search(batch_embedding_user_positive, K, exclude_links)
```

- We pass user embedding vertices  $\mathbf{h}'_{\text{user}}$  as query vectors  $Q = \{\vec{q}_0, \vec{q}_1, \dots, \vec{q}_d\}$ ,  $K$  neighbors and links to exclude.

- Recall at K calculates the fraction of relevant items among the top-K items recommended or retrieved.

$$\text{Recall@K} = \frac{\text{Number of relevant items retrieved in top-K}}{\text{Total number of relevant items}}$$

- A high recall at K indicates that the system is good at retrieving relevant items in the top-K list.

- We first initialize the class method with top-K amount of retrievals:

```
1 | recall = LinkPredRecall(k=K).to(device)
```

- Then, we pass our matrix of predicted indices from our KNN-MIPS and corresponding edge index with excluded links:

```
1 | recall.update(pred_index_mat, edge_index)
2 |
3 | recall = recall.compute()
```



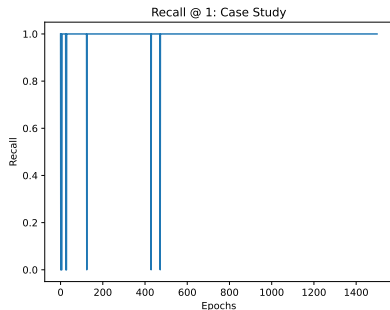
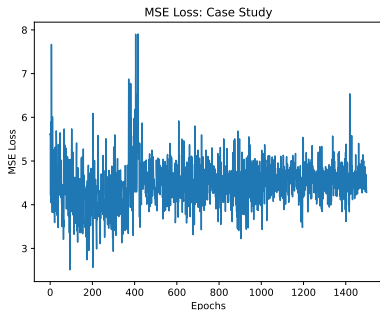
# Case Study: Simplified Movie Recommendation System

- Finally, we'll execute our GAT model with KNN-MIPS to assess whether recommending 'The Avengers' to Tim is feasible.

# Case Study: Hyperparameters

```
1 | ##### Hyperparameters #####
2 |
3 | SEED = 123
4 | SPLIT = 7/8
5 | BATCH_SIZE = 1
6 | NEIGHBORS = [1]
7 | NEG_SAMPLING = 0
8 |
9 | EPOCHS = 1500
10 | LR = 1e-02
11 | LATENT_DIM = 12
12 | HEADS = 2
13 | K = 1
14 |
15 | ##### Hyperparameters #####
```

# Case Study: MSE Loss & Recall @ 1



- Notice that after the first 500 epochs, the recommendation system converges to an optimal recommendation.



- [PyG Documentation](#)
- [PyG GAT Documentation](#)
- [PyG GAT Source Code](#)
- [KNN-MIPS Research Paper](#)
- [KNN-MIPS Lecture](#)
- [Code: Case Study - Simplified Movie Recommendation System](#)