

BackProp GCN

End-to-End Training Through Differentiable Ranking

Ian Bezerra

2025

1 The Problem

Normal GCN training uses frozen features! Normally from a vision model like ViT or Dino...

Here I will put some code snippets of the idea of making the whole process fit into one backward call and be able to run an optimizer like AdamW through the whole GCN into the Extractor!

1.1 How We Normally Do Things

```
1 # Extract features...
2 features = dinov2(images) # No gradients!
3 features = features.detach()
4
5 # Compute rankings...
6 rankings = compute_rankings(features)
7
8 # Train GCN on FIXED graph!!
9 for epoch in range(100):
10     loss = gcn(features, rankings)
11     loss.backward() # Only updates GCN, NOT DINoV2!
```

The key point here is that the Extractor doesn't learn anything. The features might not be optimal for the GCN task!

2 The Idea: Differentiable Ranking

2.1 Step 1: Compute Similarity (Differentiable)

```
1 # Normalize features
2 features_norm = F.normalize(features, p=2, dim=1)
3
4 # Cosine similarity matrix
5 similarity = features_norm @ features_norm.T
```

2.2 Step 2: Apply Soft Top-K (Differentiable)

```
1 # Apply temperature and softmax
2 soft_weights = F.softmax(similarity / temperature, dim=1)
3
4 # Get top-k indices
5 _, rankings = torch.topk(similarity, k=10)
```

```

6 # Extract soft weights for top-k neighbors
7 top_k_weights = soft_weights.gather(1, rankings)

```

We use softmax instead of hard topk.

2.3 Step 3: Build Graph

```

1 # Create edges from rankings
2 for i in range(N):
3     for j in range(k):
4         neighbor = rankings[i, j]
5         edges.append([i, neighbor])

```

The edge structure depends on rankings, which depend on features!

2.4 Step 4: GCN Forward

```

1 # Features flow through GCN
2 x = features # [N, D] - has gradients!
3 x = gcn_layer1(x, edges)
4 x = relu(x)
5 x = gcn_layer2(x, edges)
6 output = log_softmax(x)

```

Everything is differentiable.

2.5 Step 5: Backward Pass

```

1 loss = nll_loss(output[train_mask], labels[train_mask])
2 loss.backward()
3
4 # Gradients flow backward through:
5 # loss -> output -> gcn_layer2 -> gcn_layer1 -> features
6 #       -> soft_weights -> similarity -> features_norm
7 #       -> features -> dinov2!

```

3 End-to-End Training Loop

```

1 # Optimizer includes BOTH GCN and DINOV2
2 optimizer = Adam(gcn.parameters() + dinov2.parameters())
3
4 for epoch in range(200):
5     # 1. Extract features (with gradients!)
6     features = dinov2(images) # Gradients enabled!
7
8     # 2. Compute differentiable rankings
9     similarity = features_norm @ features_norm.T
10    soft_weights = softmax(similarity / temperature)
11    rankings, top_k_weights = topk(similarity, soft_weights, k=10)
12
13    # 3. Build graph from rankings
14    edges = build_edges(rankings)
15
16    # 4. Forward through GCN

```

```

17     output = gcn(features, edges)
18
19     # 5. Compute loss
20     loss = nll_loss(output[train_mask], labels)
21
22     # 6. Backward (updates BOTH GCN and DINoV2!)
23     loss.backward()
24
25     # 7. Update parameters
26     optimizer.step()  # Updates both GCN AND DINoV2!

```

4 Math: Differentiable Top-K Selection

4.1 The Problem with Hard Top-K

Traditional ranking uses discrete operations:

```

1 indices = argsort(scores) [:k]

```

This has **zero gradient everywhere**:

$$\frac{\partial \text{topk}(x)}{\partial x} = 0 \quad (1)$$

No gradients means no backpropagation!

4.2 Our Solution: Soft Top-K with Softmax

Instead, we use a **soft approximation** via temperature-scaled softmax:

4.2.1 Step 1: Compute Similarity Matrix

For features $F \in \mathbb{R}^{N \times D}$:

$$S_{ij} = \frac{F_i \cdot F_j}{\|F_i\| \|F_j\|} \quad (2)$$

This is the **cosine similarity** between feature vectors.

4.2.2 Step 2: Apply Temperature-Scaled Softmax

For each node i , compute soft weights over all nodes:

$$w_{ij} = \frac{\exp(S_{ij}/\tau)}{\sum_{l=1}^N \exp(S_{il}/\tau)} \quad (3)$$

Where:

- τ is the **temperature** parameter
- Lower $\tau \rightarrow$ sharper distribution (closer to hard selection)
- Higher $\tau \rightarrow$ smoother distribution

4.2.3 Step 3: Select Top-K Neighbors

Get indices of top-k similarities:

$$\mathcal{N}_k(i) = \text{topk}(S_i, k) \quad (4)$$

Extract corresponding soft weights:

$$\tilde{w}_{ij} = w_{ij}, \quad \forall j \in \mathcal{N}_k(i) \quad (5)$$

These weights \tilde{w}_{ij} are **differentiable!**

4.3 Why This Works: The Gradient

The softmax function has well-defined gradients:

$$\frac{\partial w_{ij}}{\partial S_{il}} = \frac{1}{\tau} w_{ij} (\delta_{jl} - w_{il}) \quad (6)$$

Where δ_{jl} is the Kronecker delta.

This means:

- **Non-zero gradients** flow backward through the softmax
- Gradients propagate from $w_{ij} \rightarrow S_{ij} \rightarrow F_i \rightarrow \text{DINOv2}$
- The entire pipeline remains differentiable!

4.4 Temperature Parameter Effect

The temperature τ controls the approximation quality:

$$\lim_{\tau \rightarrow 0} w_{ij} = \begin{cases} 1 & \text{if } j = \arg \max_l S_{il} \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

In practice:

- $\tau = 0.1$: Sharp, almost discrete (gradients can vanish)
- $\tau = 1.0$: Balanced (recommended)
- $\tau = 10.0$: Very smooth (weak selection)

4.5 Complete Forward Pass Formula

For node i with features f_i :

$$S_i = f_i \cdot F^T / (\|f_i\| \|F\|) \quad (\text{similarity}) \quad (8)$$

$$w_i = \text{softmax}(S_i / \tau) \quad (\text{soft weights}) \quad (9)$$

$$\mathcal{N}_k(i) = \text{topk}(S_i, k) \quad (\text{neighbor indices}) \quad (10)$$

$$\tilde{w}_i = w_i[\mathcal{N}_k(i)] \quad (\text{top-k weights}) \quad (11)$$

All operations are differentiable with respect to f_i !

4.6 Backward Pass (Chain Rule)

Given loss \mathcal{L} , gradients flow:

$$\frac{\partial \mathcal{L}}{\partial f_i} = \sum_{j \in \mathcal{N}_k(i)} \frac{\partial \mathcal{L}}{\partial \tilde{w}_{ij}} \frac{\partial \tilde{w}_{ij}}{\partial w_{ij}} \frac{\partial w_{ij}}{\partial S_{ij}} \frac{\partial S_{ij}}{\partial f_i} \quad (12)$$

Every term is **non-zero and computable!**