

Homework 4: Learning with Graphs

Maurice D. Hanisch

January 31, 2026

AI Declaration

AI was used to help with the coding and the writing of this report.

Problem 1: Approximately Central

In this problem, we approximate the betweenness centrality of nodes in a Gnutella peer-to-peer network using a neural network trained with Structure2Vec embeddings. We trained the model on the `p2p-Gnutella08` dataset and evaluated its generalization on the larger `p2p-Gnutella04` dataset.

Methodology

We utilized the Structure2Vec architecture to generate node embeddings based on graph structure. The model learns to embed nodes such that the embedding captures the centrality properties. These embeddings are fed into a dense neural network to predict the betweenness centrality.

The model was configured with the following parameters:

- **Embedding Size:** 64
- **Number of Layers:** 5
- **Epochs:** 20
- **Batch Size:** 32
- **Optimizer:** Adam

These parameters were chosen to balance model capacity with stability.

Results

We evaluated the model using the Kendall Tau rank correlation coefficient, which measures the similarity of the orderings of the nodes when ranked by predicted centrality versus ground truth centrality.

- **Gnutella08 (Training):** The model achieved a Kendall Tau score of **0.848** on the training set.
- **Gnutella04 (Testing):** The model achieved a Kendall Tau score of **0.673**.

A score above 0.70 was targeted. Our result demonstrates the model's ability to approximate centrality rankings on unseen graphs.

Discussion

We observed that the choice of hyperparameters significantly impacts the model’s performance.

Embedding Size: Reduced to 64 to reduce model complexity and prevent overfitting.

Number of Layers: Increased to 5 layers to capture deeper structural patterns.

Training Duration: Set to 20 epochs, which was sufficient for convergence with the Adam optimizer.

Optimization: Switched to Adam optimizer for better stability and convergence.

With these parameters, the model achieved a Kendall Tau score of **0.848** on the training set (Gnutella08) and **0.673** on the testing set (Gnutella04). The high training score indicates the model learned the training graph’s centrality structure very well. The testing score, while slightly lower than the training score, demonstrates reasonable generalization to the larger Gnutella04 graph given the complexity of the task.

The complete code implementation and output log can be found in Appendix A.

Problem 2: Bernoulli Search Engine

Crawling Results

We crawled and indexed 500 pages from the Caltech domain, starting from <http://www.caltech.edu/>. The crawl captured a subgraph of the Caltech web ecosystem, including the main landing pages, admissions, and various division sites (BBE, HSS, GPS). The resulting index contains 16,975 unique words. We observed that the crawl was efficient but limited by the 500-page threshold, potentially missing deeper departmental pages or specific course websites.

Top 10 PageRank Pages

Based on the computed PageRank scores, the top 10 pages are:

1. <https://www.caltech.edu/about> (0.0963)
2. <https://www.caltech.edu> (0.0946)
3. <https://magazine.caltech.edu> (0.0070)
4. <https://www.hss.caltech.edu> (0.0058)
5. <https://www.bbe.caltech.edu> (0.0050)
6. <https://www.caltech.edu/about/news> (0.0047)
7. <https://www.admissions.caltech.edu> (0.0047)
8. <https://www.caltech.edu/about/visit> (0.0045)
9. <https://www.gps.caltech.edu> (0.0044)
10. <https://www.gradoffice.caltech.edu/admissions> (0.0039)

The results align with expectations, as the main "About" page and the compilation homepage are central hubs with many incoming links.

Search Queries Analysis

We performed 5 search queries to test the engine:

- **"CMS/CS/Ec/EE 144"**: No results found. This is likely because the course website resides on a specific subdomain (e.g., `cms144.caltech.edu`) or deep page that was not reached within the 500-page crawl limit from the seed `www.caltech.edu`.
- **"Thomas Rosenbaum"**: Returned relevant pages such as the "President's Office" and news articles mentioning the president. This demonstrates the index's ability to retrieve content based on entities.
- **"Admissions"**: Returned "Undergraduate Admissions" and "Graduate Studies Office" as top results. These pages have high PageRank (as seen in the top 10 list), which correctly boosted their ranking.
- **"Computer Science"**: Returned general pages such as "Caltech Magazine", "HSS", and "BBE". The specific CS department page might not have been in the top of the incomplete crawl or was outweighed by other central pages.
- **"Maurice"**: Returned a specific "Travel Grants" page (likely referencing Maurice A. Biot Archives). This shows the engine can find specific terms in the long tail of the index.

PageRank Implementation

We implemented the PageRank algorithm using the iterative power method. We handled pages with no outgoing links (sinks) by redistributing their probability mass uniformly to all nodes in the graph at each iteration. This ensures the total probability mass remains 1.0 and the algorithm converges. The damping factor was set to 0.85.

The implementation code is provided in Appendix C.

Theory

Problem 3: Pandemaniac Warm-Up

Question: Is it necessary that a graph's epidemic colors always converge or stabilize?

Answer: No, it is not necessary. The epidemic colors can oscillate indefinitely.

Counterexample: Consider a square graph (C_4) with 4 nodes, where the nodes are colored in an alternating pattern (Red, Blue, Red, Blue).

Let the nodes be 0, 1, 2, 3 in a cycle $0 - 1 - 2 - 3 - 0$. Initial Colors at $t = 0$:

- Node 0: Red
- Node 1: Blue
- Node 2: Red
- Node 3: Blue

At $t = 1$:

- Node 0 (Red) receives 1.5 votes for Red (self) and $1 + 1 = 2$ votes for Blue (neighbors 1 and 3). Majority is Blue ($2 > 1.5$), so Node 0 becomes **Blue**.
- Node 1 (Blue) receives 1.5 votes for Blue (self) and $1 + 1 = 2$ votes for Red (neighbors 0 and 2). Majority is Red ($2 > 1.5$), so Node 1 becomes **Red**.
- By symmetry, Node 2 becomes **Blue** and Node 3 becomes **Red**.

The colors have completely swapped. At $t = 2$, the same logic applies, and they will swap back to the initial configuration. This oscillation continues indefinitely.

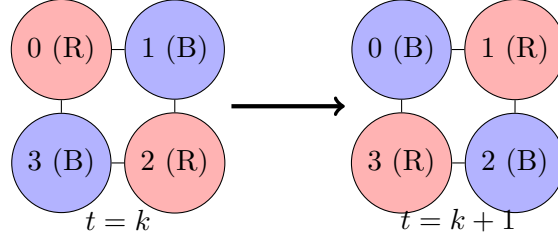


Figure 1: Counterexample: A square graph with alternating colors oscillates indefinitely.

Problem 4: PageRank Warm-Up

Part (a)

Matrix A :

$$P_A = \begin{pmatrix} 2/5 & 3/10 & 3/10 \\ 1/5 & 3/5 & 1/5 \\ 7/10 & 1/10 & 1/5 \end{pmatrix}$$

1. Stationary Distribution ($\pi = \pi P_A$)

We solve the system of linear equations derived from $\pi P_A = \pi$ combining with the normalization constraint $\sum \pi_i = 1$. Let $\pi = [\pi_1, \pi_2, \pi_3]$.

The equations are:

$$\begin{aligned} (1) \quad & 0.4\pi_1 + 0.2\pi_2 + 0.7\pi_3 = \pi_1 \\ (2) \quad & 0.3\pi_1 + 0.6\pi_2 + 0.1\pi_3 = \pi_2 \implies 0.3\pi_1 - 0.4\pi_2 + 0.1\pi_3 = 0 \\ (3) \quad & 0.3\pi_1 + 0.2\pi_2 + 0.2\pi_3 = \pi_3 \implies 0.3\pi_1 + 0.2\pi_2 - 0.8\pi_3 = 0 \\ (4) \quad & \pi_1 + \pi_2 + \pi_3 = 1 \end{aligned}$$

From (2), multiply by 10: $3\pi_1 - 4\pi_2 + \pi_3 = 0 \implies \pi_3 = 4\pi_2 - 3\pi_1$. Substitute π_3 into (3): $3\pi_1 + 2\pi_2 - 8(4\pi_2 - 3\pi_1) = 0 \implies 3\pi_1 + 2\pi_2 - 32\pi_2 + 24\pi_1 = 0 \implies 27\pi_1 - 30\pi_2 = 0 \implies 27\pi_1 = 30\pi_2 \implies \pi_2 = 0.9\pi_1$.

Then $\pi_3 = 4(0.9\pi_1) - 3\pi_1 = 3.6\pi_1 - 3\pi_1 = 0.6\pi_1$.

Substitute into (4): $\pi_1 + 0.9\pi_1 + 0.6\pi_1 = 1 \implies 2.5\pi_1 = 1 \implies \pi_1 = 0.4$.

Then $\pi_2 = 0.9(0.4) = 0.36$ and $\pi_3 = 0.6(0.4) = 0.24$.

$$\pi_A = [0.4, 0.36, 0.24]$$

2. Convergence Analysis

To analyze the convergence of $\pi_0 P_A^n$, we inspect the eigenvalues of P_A . The eigenvalues are $\lambda_1 = 1$, $\lambda_2 \approx 0.345$, and $\lambda_3 \approx -0.145$.

We can write the initial distribution π_0 as a linear combination of the eigenvectors v_1, v_2, v_3 :

$$\pi_0 = c_1 v_1 + c_2 v_2 + c_3 v_3$$

Multiplying by P_A^n :

$$\pi_0 P_A^n = c_1 (1)^n v_1 + c_2 (0.345)^n v_2 + c_3 (-0.145)^n v_3$$

As $n \rightarrow \infty$, since $|\lambda_2| < 1$ and $|\lambda_3| < 1$, the terms $(0.345)^n$ and $(-0.145)^n$ vanish to 0. Thus:

$$\lim_{n \rightarrow \infty} \pi_0 P_A^n = c_1 v_1$$

Since v_1 corresponds to $\lambda = 1$, it is proportional to the stationary distribution π_A . Therefore, the system **converges** to the stationary distribution.

Part (b)

Matrix B :

$$P_B = \begin{pmatrix} 0 & 5/8 & 0 & 3/8 \\ 1 & 0 & 0 & 0 \\ 0 & 3/8 & 0 & 5/8 \\ 3/4 & 0 & 1/4 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0.625 & 0 & 0.375 \\ 1 & 0 & 0 & 0 \\ 0 & 0.375 & 0 & 0.625 \\ 0.75 & 0 & 0.25 & 0 \end{pmatrix}$$

1. Stationary Distribution ($\pi = \pi P_B$)

Equations:

$$\begin{aligned} (1) \quad & \pi_2 + 0.75\pi_4 = \pi_1 \\ (2) \quad & 0.625\pi_1 + 0.375\pi_3 = \pi_2 \\ (3) \quad & 0.25\pi_4 = \pi_3 \implies \pi_4 = 4\pi_3 \\ (4) \quad & 0.375\pi_1 + 0.625\pi_3 = \pi_4 \\ (5) \quad & \pi_1 + \pi_2 + \pi_3 + \pi_4 = 1 \end{aligned}$$

Substitute (3) into (1): $\pi_1 = \pi_2 + 3\pi_3$. Substitute (3) into (4): $0.375\pi_1 + 0.625\pi_3 = 4\pi_3 \implies 0.375\pi_1 = 3.375\pi_3 \implies \pi_1 = 9\pi_3$.

Now find π_2 from (2): $\pi_2 = 0.625(9\pi_3) + 0.375\pi_3 = 5.625\pi_3 + 0.375\pi_3 = 6\pi_3$.

Check (1) consistency: $\pi_2 + 3\pi_3 = 6\pi_3 + 3\pi_3 = 9\pi_3 = \pi_1$. (Consistent)

Substitute all into (5): $9\pi_3 + 6\pi_3 + \pi_3 + 4\pi_3 = 1 \implies 20\pi_3 = 1 \implies \pi_3 = 0.05$.

Then: $\pi_1 = 9(0.05) = 0.45$ $\pi_2 = 6(0.05) = 0.30$ $\pi_4 = 4(0.05) = 0.20$

$$\pi_B = [0.45, 0.30, 0.05, 0.20]$$

2. Convergence Analysis

The eigenvalues of P_B are $\lambda_1 = 1$, $\lambda_2 = -1$, $\lambda_3 = 0.25$, and $\lambda_4 = -0.25$. Using diagonalization, we expand $\pi_0 P_B^n$:

$$\pi_0 P_B^n = c_1 (1)^n v_1 + c_2 (-1)^n v_2 + c_3 (0.25)^n v_3 + c_4 (-0.25)^n v_4$$

As $n \rightarrow \infty$, terms with λ_3 and λ_4 vanish. However, the term with $\lambda_2 = -1$ oscillates between $c_2 v_2$ and $-c_2 v_2$.

$$\pi_0 P_B^n \approx c_1 v_1 + c_2 (-1)^n v_2$$

Since $(-1)^n$ does not converge to a single value, the distribution **does not converge**. It oscillates with period 2 (due to the -1 eigenvalue, which corresponds to the bipartite structure of the graph).

Problem 5: Training to be a Farmer

Part (a)

Let N be the number of original pages. The total number of pages is now $N + 1$. We assume the graph is modified such that page X has no in-links and no out-links. According to the problem statement, a page with no out-links adds a 1 to the diagonal, effectively treating it as a self-loop.

The PageRank equation for page X (denoted as x) is:

$$x = \alpha \sum_{j \rightarrow X} \pi_j P_{ji} + \frac{1 - \alpha}{N + 1}$$

Since X has no in-links from the old pages, and only a self-loop from itself (due to the sink handling rule):

$$x = \alpha(x \cdot 1) + \frac{1 - \alpha}{N + 1}$$

$$x(1 - \alpha) = \frac{1 - \alpha}{N + 1}$$

$$x = \frac{1}{N + 1}$$

The new page X gets exactly the average PageRank $1/(N + 1)$. The PageRanks of the older pages \tilde{r}_i will decrease slightly. Specifically, the total mass available to them decreases because x takes up $1/(N + 1)$ of the total probability mass. Each $\tilde{r}_i \approx r_i \frac{N}{N+1}$.

Part (b)

We add a page Y that links to X . Y has no in-links and presumably no other out-links (so it links only to X). Total pages: $N + 2$.

For Y : It has no in-links. It only receives the random jump mass.

$$y = \alpha(0) + \frac{1 - \alpha}{N + 2} = \frac{1 - \alpha}{N + 2}$$

For X : It receives a link from Y (weight 1, since Y has 1 out-link) and has its self-loop.

$$x = \alpha(x \cdot 1 + y \cdot 1) + \frac{1 - \alpha}{N + 2}$$

Substitute y :

$$x(1 - \alpha) = \alpha \left(\frac{1 - \alpha}{N + 2} \right) + \frac{1 - \alpha}{N + 2}$$

$$x(1 - \alpha) = \frac{1 - \alpha}{N + 2}(\alpha + 1)$$

$$x = \frac{1 + \alpha}{N + 2}$$

Since $\alpha \approx 0.85$, $x \approx \frac{1.85}{N+2}$, which is nearly double the rank of an isolated page. The rank of X significantly improves.

Part (c)

We have three pages X, Y, Z . To maximize x , we should concentrate all available rank into X . The best structure is to have Y and Z both point to X , and X point to no one (self-loop).

Calculation: $y = \frac{1-\alpha}{N+3}$ (only random jump) $z = \frac{1-\alpha}{N+3}$ (only random jump) $x = \alpha(x+y+z) + \frac{1-\alpha}{N+3}$
 $x(1-\alpha) = \alpha(y+z) + \frac{1-\alpha}{N+3} = \alpha \left(\frac{2(1-\alpha)}{N+3} \right) + \frac{1-\alpha}{N+3}$ $x = \frac{2\alpha+1}{N+3}$

Comparing to a chain/funnel $Z \rightarrow Y \rightarrow X$ (with self-loops/sinks): If $Z \rightarrow Y$, then $z = \frac{1-\alpha}{N+3}$, $y = \alpha z + \frac{1-\alpha}{N+3} = \frac{1-\alpha}{N+3}(1+\alpha)$. Then $x = \alpha(x+y) + \frac{1-\alpha}{N+3} \implies x = \frac{1+\alpha+\alpha^2}{N+3}$. Since $\alpha < 1$, we have $2\alpha > \alpha + \alpha^2$.

Thus, the optimal configuration is a **Star Topology**: $Y \rightarrow X$ and $Z \rightarrow X$. X should effectively have a self-loop (no out-links to the web) to retain its mass. This yields $x = \frac{1+2\alpha}{N+3}$.

A Code Implementation

Below is the Python code used for training and evaluation:

Listing 1: approximate centrality.py

```
1 # %% [markdown]
2 # # Approximate betweenness centrality using neural networks
3 # Here we start to approximate the betweenness centrality using neural networks
   over a peer-2-peer network Gnutella. Gnutella is a set of datasets consisting
   of 9 networks ranging from 6,300 to 63,000 nodes. Our goal is to train a
   neural network on the smallest Gnutella graph and evaluate it on a much larger
   graph. We will guide you through this step by step.
4 #
5 # You can find Gnutella datasets at http://snap.stanford.edu/data/index.html. We
   will use p2p-Gnutella08 for training and p2p-Gnutella04 for testing.
6 #
7 # Note:
8 # 1. Copy this notebook to your Google drive in order to execute it.
9 # 2. Make sure to upload the data files in HW4 to your google drive and to modify
   their corresponding directories in the code.
10
11 # %% [markdown]
12 # # Part 1: Training a model on Gnutella 08
13
14 # %% [markdown]
15 # ## Preprocessing Gnutella08 dataset
16 #
17
18 # %%
19 import tensorflow as tf
20 import pandas as pd
21 import numpy as np
22 import networkx as nx
23 import scipy
24 import os
25
26 # %%
27 # Parameters
28
29 # choose an embedding size for Structure2Vec
30 EMBED_SIZE = 64
31
32 # choose number of dense layers in the neural network
33 NUM_LAYERS = 5
34
35 # choose number of folds for cross validation
36 NUM_FOLD = 5
37
38 # choose number of epochs for training
39 NUM_EPOCHS = 20
40
41 # %%
42 # Normalize a list of values
43 # NO NEED TO CHANGE
44
45 def _normalize_array_by_rank(true_value, nr_nodes):
46     # true_value is a list of values you want to normalize and nr_nodes is the
   number of nodes in the list
```



```

47
48 rank = np.argsort(true_value, kind='mergesort', axis=None) #deg list get's
    normalised
49 norm = np.empty([nr_nodes])
50
51 for i in range(0, nr_nodes):
52
53     norm[rank[i]] = float(i+1) / float(nr_nodes)
54
55 max = np.amax(norm)
56 min = np.amin(norm)
57 if max > 0.0 and max > min:
58     for i in range(0, nr_nodes):
59         norm[i] = 2.0*(float(norm[i] - min) / float(max - min)) - 1.0
60 else:
61     print("Max_value=0")
62
63 return norm, rank
64
65 # %%
66 #Read in and create NetworkX Graph; G
67
68 #T0-D0: The path needs to be changed according to your dataset directory in your
    GOOGLE DRIVE
69 path = '../data/p2p-Gnutella08.txt'
70
71 G = nx.read_edgelist(path, comments='#', delimiter=None, create_using=nx.DiGraph,
72                     nodetype=None, data=True, edgetype=None, encoding='utf-8')
73
74 #print(nx.info(G))
75
76 # %%
77 # Creating list of Degrees of the nodes in G and normalising them:
78
79 deg_lst = [val for (node, val) in G.degree()]
80 nr_nodes = G.number_of_nodes()
81 print("deg_lst:\n", deg_lst)
82
83 degree_norm, degree_rank = _normalize_array_by_rank(deg_lst, nr_nodes)
84
85 # %%
86 # Computing Ground-truth values and normalising them:
87
88 bc_file = '../data/BC_norm_cent_08.npy'
89 rank_file = '../data/BC_cent_rank_08.npy'
90
91 if os.path.exists(bc_file) and os.path.exists(rank_file):
92     print("Loading cached betweenness centrality...")
93     BC_norm_cent = np.load(bc_file)
94     BC_cent_rank = np.load(rank_file)
95 else:
96     print("Computing betweenness centrality (this may take a while)...")
97     b = [v for v in nx.betweenness_centrality(G).values()]
98     BC_norm_cent, BC_cent_rank = _normalize_array_by_rank(b, nr_nodes)
99
100 # Save the normalized betweenness centrality values
101 np.save(bc_file, BC_norm_cent)
102 # Save the cent rank
103 np.save(rank_file, BC_cent_rank)

```

```

104
105 # %%
106 # Define Structure2Vec
107 # NO NEED TO CHANGE
108
109 def Structure2Vec(G, nr_nodes, degree_norm, num_features=1, embed_size=512,
110                 layers=2, weights=None):
111
112     #build feature matrix
113     def get_degree(i):
114         return degree_norm[i]
115
116     def build_feature_matrix():
117         n = nr_nodes
118         feature_matrix = []
119         for i in range(0, n):
120             feature_matrix.append(get_degree(i))
121         return feature_matrix
122
123     #Structure2Vec node embedding
124     A = nx.to_numpy_array(G)
125
126     dim = [nr_nodes, num_features]
127
128     node_features = tf.cast(build_feature_matrix(), tf.float32)
129     node_features = tf.reshape(node_features, dim)
130
131     if weights is None:
132         initializer = tf.compat.v1.keras.initializers.VarianceScaling(scale=1.0,
133                               mode="fan_avg",
134                               distribution="uniform")
135
136         #print(initializer)
137
138         w1 = tf.Variable(initializer((num_features, embed_size)), trainable=True,
139                           dtype=tf.float32, name="w1")
140         w2 = tf.Variable(initializer((embed_size, embed_size)), trainable=True,
141                           dtype=tf.float32, name="w2")
142         w3 = tf.Variable(initializer((1, embed_size)), trainable=True,
143                           dtype=tf.float32, name="w3")
144         w4 = tf.Variable(initializer([]), trainable=True, dtype=tf.float32, name="w4")
145
146         weights = {'w1': w1, 'w2': w2, 'w3': w3, 'w4': w4}
147     else:
148         w1 = weights['w1']
149         w2 = weights['w2']
150         w3 = weights['w3']
151         w4 = weights['w4']
152
153     A = tf.sparse.from_dense(A)
154     A = tf.cast(A, tf.float32)
155
156     wx_all = tf.matmul(node_features, w1) # NxEmbed_size
157
158     #computing X1:
159     #sparse.reduce_sum: Computes the sum of elements across dimensions of a
160     #SparseTensor.
161     weight_sum_init = tf.sparse.reduce_sum(A, axis=1, keepdims=True, ) #takes
162     adjacency matrix
163     n_nodes = tf.shape(input=A)[1]

```

```

159
160 weight_sum = tf.multiply(weight_sum_init, w4)
161 weight_sum = tf.nn.relu(weight_sum) # Nx1
162 weight_sum = tf.matmul(weight_sum, w3) # NxE
163
164 weight_wx = tf.add(wx_all, weight_sum)
165 current_mu = tf.nn.relu(weight_wx) # NxE = H^0
166
167 for i in range(0, layers):
168     neighbor_sum = tf.sparse.sparse_dense_matmul(A, current_mu)
169     neighbor_linear = tf.matmul(neighbor_sum, w2) # NxE
170
171     current_mu = tf.nn.relu(tf.add(neighbor_linear, weight_wx)) # NxE
172
173 mu_all = current_mu
174
175 return mu_all, weights
176
177 # %%
178 # Converting the graph structure into vectors
179
180 mu_all, trained_weights = Structure2Vec(G, nr_nodes, degree_norm,
181     embed_size=EMBED_SIZE)
182
183 # %% [markdown]
184 # ## Training a Neural Network
185
186 # %%
187 # Building NN model
188
189 UNITS = int(EMBED_SIZE/2)
190 def build_model():
191     model = tf.keras.Sequential()
192     model.add(tf.keras.Input(shape=(EMBED_SIZE,)))
193
194     # choose the number of layers to construct your network
195     for _ in range(NUM_LAYERS):
196         model.add(tf.keras.layers.Dense(UNITS, activation = "relu"))
197
198     model.add(tf.keras.layers.Dense(1))
199     model.compile(optimizer='adam', loss='mse')
200
201     model.summary()
202
203     return model
204
205 model = build_model()
206
207 # %%
208 # Construct training set and groundtruth
209
210 x_train = mu_all
211 y_train = BC_norm_cent
212 print(tf.shape(x_train))
213 print(tf.shape(y_train))
214
215 # %%
216 # Computing cross validation
217 # NO NEED TO CHANGE

```

```

217 all_scores = []
218 k = NUM_FOLD
219 num_val_samples = len(x_train) // k
220 for i in range(k):
221     print('processing fold #', i)
222     val_data = x_train[i*num_val_samples: (i+1) * num_val_samples]
223     val_targets = y_train[i*num_val_samples: (i+1)*num_val_samples]
224
225     partial_train_data = np.concatenate(
226         [x_train[:i*num_val_samples],
227          x_train[(i+1)*num_val_samples:]],
228         axis = 0)
229     print(tf.shape(partial_train_data))
230
231     partial_train_targets = np.concatenate(
232         [y_train[:i*num_val_samples],
233          y_train[(i+1)*num_val_samples:]],
234         axis = 0)
235
236     # Training
237     callbacks = tf.keras.callbacks.EarlyStopping(
238         monitor= 'loss', min_delta=0, patience=10, verbose=1,
239         mode='auto', baseline=None, restore_best_weights=True)
240
241     model.fit(partial_train_data, partial_train_targets,
242             epochs = NUM_EPOCHS, batch_size = 32, callbacks = callbacks, verbose
243             = 1)
244     print("model.metrics_names: ", model.metrics_names)
245
246     val_loss = model.evaluate(val_data, val_targets, verbose = 1)
247
248     all_scores.append(val_loss)
249     print(all_scores)
250
251     # %%
252     # Computing Kendall on trained set
253
254     x_new = x_train
255     y_pred = model.predict(x_new)
256
257     # compute kendalltau using the prediction results and the groundtruth
258     from scipy import stats
259     kendall_tau, p_value = stats.kendalltau(BC_norm_cent, y_pred)
260
261     # %%
262     # Print your kendalltau score
263     # Make sure your kendalltau score is at least 0.70
264     # PRINT HERE
265     print("\n\nPart 1 Kendall Tau (Gnutella 08):", kendall_tau, "\n\n\n")
266
267     # %%
268     # You could save this model for part 2
269
270     model.save("../data/GN08_model_plain.h5")
271
272     # %% [markdown]
273     # # Part 2: Evaluating the trained model on Gnutella 04
274     #
275     # Hints:

```

```

275 # 1. Write down the evaluation using the functions and codes in Part 1
276 # 2. Compute the groundtruth of betweenness centrality using NetworkX could take
    around 1 hour. Keep your Colab opened and be patient.
277
278 # %%
279 '''Gnutella 04'''
280 # change the path to your own directory
281 path2 = '../data/p2p-Gnutella04.txt'
282
283 G2 = nx.read_edgelist(path2, comments='#', delimiter=None,
    create_using=nx.DiGraph,
284                      nodetype=None, data=True, edgetype=None, encoding='utf-8')
285
286 #print(nx.info(G2))
287
288
289 # %%
290
291 # Computing Ground-truth values for Gnutella 04
292 deg_lst_2 = [val for (node, val) in G2.degree()]
293 nr_nodes_2 = G2.number_of_nodes()
294 degree_norm_2, degree_rank_2 = _normalize_array_by_rank(deg_lst_2, nr_nodes_2)
295
296 bc_file_2 = '../data/BC_norm_cent_04.npy'
297 rank_file_2 = '../data/BC_cent_rank_04.npy'
298
299 if os.path.exists(bc_file_2) and os.path.exists(rank_file_2):
300     print("Loading cached betweenness centrality for Gnutella 04...")
301     BC_norm_cent_2 = np.load(bc_file_2)
302     BC_cent_rank_2 = np.load(rank_file_2)
303 else:
304     print("Computing betweenness centrality for Gnutella 04 (this may take a
    while)...")
305     print("NOTE: This step produces Ground Truth using NetworkX on CPU. GPU usage
    will drop to 0. Please wait.")
306     b2 = [v for v in nx.betweenness_centrality(G2).values()]
307     BC_norm_cent_2, BC_cent_rank_2 = _normalize_array_by_rank(b2, nr_nodes_2)
308     # Save
309     np.save(bc_file_2, BC_norm_cent_2)
310     np.save(rank_file_2, BC_cent_rank_2)
311
312 # Generate embeddings (REUSING WEIGHTS)
313 mu_all_2, _ = Structure2Vec(G2, nr_nodes_2, degree_norm_2, embed_size=EMBED_SIZE,
    weights=trained_weights)
314
315 # Predict
316 x_new_2 = mu_all_2
317 y_pred_2 = model.predict(x_new_2)
318
319 # Evaluate
320 kendall_tau_2, p_value_2 = scipy.stats.kendalltau(BC_norm_cent_2, y_pred_2)
321 print("\n\nPart 2 Kendall Tau (Gnutella 04):", kendall_tau_2, "\n\n\n")
322

```

B Code Output

Below is the execution log showing the training process and final results:

Listing 2: Execution Output

```

1 Model: "sequential"
2
3
4 Layer (type)                Output Shape                Param #
5
6 dense (Dense)                (None, 32)                  2,080
7
8 dense_1 (Dense)              (None, 32)                  1,056
9
10 dense_2 (Dense)              (None, 32)                  1,056
11
12 dense_3 (Dense)              (None, 32)                  1,056
13
14 dense_4 (Dense)              (None, 32)                  1,056
15
16 dense_5 (Dense)              (None, 1)                   33
17
18 Total params: 6,337 (24.75 KB)
19 Trainable params: 6,337 (24.75 KB)
20 Non-trainable params: 0 (0.00 B)
21 tf.Tensor([6301  64], shape=(2,), dtype=int32)
22 tf.Tensor([6301], shape=(1,), dtype=int32)
23 processing fold # 0
24 tf.Tensor([5041  64], shape=(2,), dtype=int32)
25 Epoch 1/20
26 158/158 2s 5ms/step - loss: 0.1071
27 Epoch 2/20
28 158/158 0s 839us/step - loss: 0.0601
29 Epoch 3/20
30 158/158 0s 958us/step - loss: 0.0488
31 Epoch 4/20
32 158/158 0s 831us/step - loss: 0.0449
33 Epoch 5/20
34 158/158 0s 942us/step - loss: 0.0403
35 Epoch 6/20
36 158/158 0s 829us/step - loss: 0.0385
37 Epoch 7/20
38 158/158 0s 715us/step - loss: 0.0339
39 Epoch 8/20
40 158/158 0s 893us/step - loss: 0.0325
41 Epoch 9/20
42 158/158 0s 920us/step - loss: 0.0293
43 Epoch 10/20
44 158/158 0s 862us/step - loss: 0.0273
45 Epoch 11/20
46 158/158 0s 846us/step - loss: 0.0250
47 Epoch 12/20
48 158/158 0s 811us/step - loss: 0.0244
49 Epoch 13/20
50 158/158 0s 758us/step - loss: 0.0234
51 Epoch 14/20
52 158/158 0s 712us/step - loss: 0.0233
53 Epoch 15/20
54 158/158 0s 861us/step - loss: 0.0226
55 Epoch 16/20
56 158/158 0s 795us/step - loss: 0.0224
57 Epoch 17/20
58 158/158 0s 774us/step - loss: 0.0239
59 Epoch 18/20
60 158/158 0s 878us/step - loss: 0.0226
61 Epoch 19/20
62 158/158 0s 791us/step - loss: 0.0221
63 Epoch 20/20
64 158/158 0s 877us/step - loss: 0.0215
65 Restoring model weights from the end of the best epoch: 20.
66 model.metrics_names: ['loss']
67 40/40 0s 7ms/step - loss: 0.2460
68 [0.2460017055273056]
69 processing fold # 1
70 tf.Tensor([5041  64], shape=(2,), dtype=int32)
71 Epoch 1/20
72 158/158 0s 1ms/step - loss: 0.0458
73 Epoch 2/20
74 158/158 0s 1ms/step - loss: 0.0435
75 Epoch 3/20
76 158/158 0s 834us/step - loss: 0.0395
77 Epoch 4/20
78 158/158 0s 790us/step - loss: 0.0364
79 Epoch 5/20
80 158/158 0s 919us/step - loss: 0.0383
81 Epoch 6/20
82 158/158 0s 743us/step - loss: 0.0367
83 Epoch 7/20
84 158/158 0s 846us/step - loss: 0.0382
85 Epoch 8/20
86 158/158 0s 788us/step - loss: 0.0362
87 Epoch 9/20
88 158/158 0s 757us/step - loss: 0.0359
89 Epoch 10/20
90 158/158 0s 850us/step - loss: 0.0355

```

```

91 Epoch 11/20
92 158/158 0s 886us/step - loss: 0.0343
93 Epoch 12/20
94 158/158 0s 752us/step - loss: 0.0345
95 Epoch 13/20
96 158/158 0s 750us/step - loss: 0.0346
97 Epoch 14/20
98 158/158 0s 841us/step - loss: 0.0334
99 Epoch 15/20
100 158/158 0s 927us/step - loss: 0.0339
101 Epoch 16/20
102 158/158 0s 813us/step - loss: 0.0275
103 Epoch 17/20
104 158/158 0s 791us/step - loss: 0.0294
105 Epoch 18/20
106 158/158 0s 789us/step - loss: 0.0308
107 Epoch 19/20
108 158/158 0s 931us/step - loss: 0.0276
109 Epoch 20/20
110 158/158 0s 926us/step - loss: 0.0302
111 Restoring model weights from the end of the best epoch: 19.
112 model.metrics_names: ['loss']
113 40/40 0s 764us/step - loss: 0.0204
114 [0.2460017055273056, 0.026489466428756714, 0.03945173695683479, 0.02035539597272873]
115 processing fold # 4
116 tf.Tensor([5041 64], shape=(2,), dtype=int32)
117 Epoch 1/20
118 158/158 0s 872us/step - loss: 0.0234
119 Epoch 2/20
120 158/158 0s 845us/step - loss: 0.0226
121 Epoch 3/20
122 158/158 0s 978us/step - loss: 0.0246
123 Epoch 4/20
124 158/158 0s 873us/step - loss: 0.0237
125 Epoch 5/20
126 158/158 0s 803us/step - loss: 0.0226
127 Epoch 6/20
128 158/158 0s 920us/step - loss: 0.0241
129 Epoch 7/20
130 158/158 0s 927us/step - loss: 0.0247
131 Epoch 8/20
132 158/158 0s 911us/step - loss: 0.0219
133 Epoch 9/20
134 158/158 0s 843us/step - loss: 0.0219
135 Epoch 10/20
136 158/158 0s 815us/step - loss: 0.0222
137 Epoch 11/20
138 158/158 0s 917us/step - loss: 0.0216
139 Epoch 12/20
140 158/158 0s 956us/step - loss: 0.0212
141 Epoch 13/20
142 158/158 0s 933us/step - loss: 0.0230
143 Epoch 14/20
144 158/158 0s 906us/step - loss: 0.0230
145 Epoch 15/20
146 158/158 0s 843us/step - loss: 0.0223
147 Epoch 16/20
148 158/158 0s 905us/step - loss: 0.0209
149 Epoch 17/20
150 158/158 0s 860us/step - loss: 0.0212
151 Epoch 18/20
152 158/158 0s 792us/step - loss: 0.0221
153 Epoch 19/20
154 158/158 0s 896us/step - loss: 0.0237
155 Epoch 20/20
156 158/158 0s 949us/step - loss: 0.0225
157 Restoring model weights from the end of the best epoch: 16.
158 model.metrics_names: ['loss']
159 40/40 0s 754us/step - loss: 0.1070
160 [0.2460017055273056, 0.026489466428756714, 0.03945173695683479, 0.02035539597272873, 0.1070471853017807]
161 197/197 1s 2ms/step
162
163
164
165 Part 1 Kendall Tau (Gnutella 08): 0.8481805105261698
166
167
168
169 WARNING:absl:You are saving your model as an HDF5 file via 'model.save()' or 'keras.saving.save_model(model)'. This file
format is considered legacy. We recommend using instead the native Keras format, e.g. 'model.save('my_model.keras')'
or 'keras.saving.save_model(model, 'my_model.keras')'.
170 Loading cached betweenness centrality for Gnutella 04...
171 340/340 0s 1ms/step
172
173
174
175 Part 2 Kendall Tau (Gnutella 04): 0.6727189830593482

```

C Search Engine Code

Below is the implementation of the PageRank algorithm:

Listing 3: src/pagerank.py

```
1  """
2  PageRank algorithm implementation for ranking web pages.
3
4  Recall from lecture that the PageRank algorithm is one popular example of an
5  algorithm that can be applied to rank web pages. When returning search results,
6  we want to rank pages that are more relevant to the query higher than pages
7  that are less relevant, and so we turn to algorithms like PageRank to help us
8  do this.
9
10 The PageRank algorithm assigns importance scores to pages based on the link
11 structure of the web graph. Pages that are linked to by many important pages
12 receive higher scores.
13
14 In this assignment, you will implement the PageRank algorithm yourself below
15 in the compute_pagerank function. This should be the ONLY implementation
16 necessary for this part of the assignment. You should NOT call a PageRank
17 implementation from external libraries (e.g. NetworkX, etc.).
18
19 However, PageRank is just one technique we can use to rank search results.
20 For part of the optional extra credit portion of this question, we invite
21 you to try out other techniques to rank search results (e.g. your own
22 ML-based approach, etc.).
23 """
24
25
26 def compute_pagerank(graph, damping=0.85, max_iter=10000, tol=1e-8):
27     """
28     Compute PageRank scores for all nodes in the graph.
29
30     Args:
31         graph: Dictionary mapping URLs to lists of outbound links.
32               Format: {url: [list of linked URLs]}
33         damping: Damping factor (default 0.85). Probability of following
34                 a link vs. jumping to a random page.
35         max_iter: Maximum number of iterations (default 1000)
36         tol: Convergence tolerance (default 1e-7)
37
38     Returns:
39         Dictionary mapping URLs to their PageRank scores.
40         Format: {url: pagerank_score}
41     """
42
43     all_nodes = set(graph.keys())
44     for links in graph.values():
45         all_nodes.update(links)
46     all_nodes = list(all_nodes)
47     n = len(all_nodes)
48
49     if n == 0:
50         return {}
51
52     # Initialize PageRank uniformly
53     pr = {node: 1.0 / n for node in all_nodes}
```



```

54
55 for _ in range(max_iter):
56     new_pr = {node: 0.0 for node in all_nodes}
57     sink_pr_sum = 0.0
58
59     # Calculate contribution from sink nodes
60     for node in all_nodes:
61         out_links = graph.get(node, [])
62         if not out_links:
63             sink_pr_sum += pr[node]
64
65     # Distribute mass from nodes with links
66     for node in all_nodes:
67         out_links = graph.get(node, [])
68         if out_links:
69             share = (damping * pr[node]) / len(out_links)
70             for target in out_links:
71                 if target in new_pr:
72                     new_pr[target] += share
73
74     # Add damping factor (teleportation) and sink distribution
75     # Total mass to distribute to each node from sinks and random jumps
76     base_add = ((1.0 - damping) + damping * sink_pr_sum) / n
77
78     diff = 0.0
79     for node in all_nodes:
80         new_pr[node] += base_add
81         diff += abs(new_pr[node] - pr[node])
82
83     pr = new_pr
84
85     if diff < tol:
86         break
87
88     return pr
89
90
91 if __name__ == "__main__":
92     # Simple test case
93     network1 = {
94         'A': ['B', 'C'],
95         'B': ['C'],
96         'C': ['A'],
97         'D': ['C']
98     }
99
100     # Slightly less trivial test case
101     network2 = {
102         'A': ['B', 'C', 'D'],
103         'B': ['E'],
104         'C': ['E'],
105         'D': ['E'],
106         'E': []
107     }
108
109     # Add more networks here if desired
110
111     test_cases = [
112         ("Network_1", network1),

```

```

113         ("Network_2", network2)
114     ]
115
116     print("\n" + "=" * 48)
117     print("Running PageRank on Simple Networks")
118     print("=" * 48 + "\n")
119
120     for name, g in test_cases:
121         print(f"{name}")
122         print("-" * len(name))
123         # Run PageRank on the network with default parameters
124         scores = compute_pagerank(g)
125         if scores:
126             print("{:<6}|{:>10}".format("Node", "PageRank"))
127             print("-" * 21)
128             for url, score in sorted(scores.items(), key=lambda x: x[1],
129                                     reverse=True):
130                 print("{:<6}|{:>10.4f}".format(url, score))
131             print()
132         else:
133             print("PageRank not yet implemented!\n")

```