

Shahd Mohamed Abd El Sabour

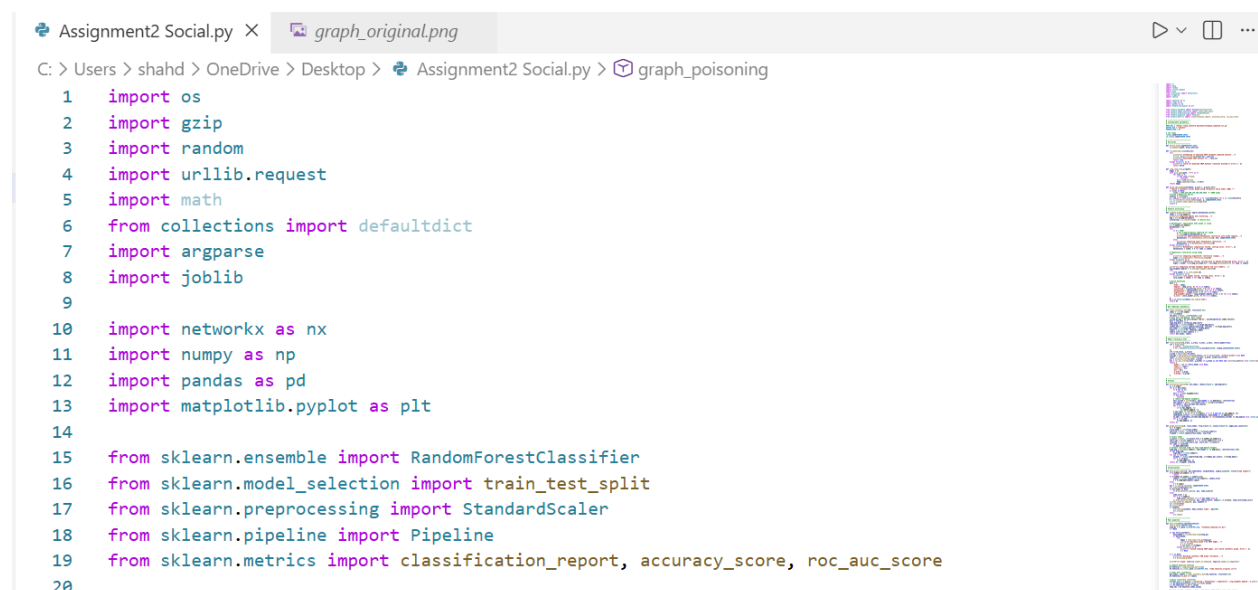
2205027

Assignment 2 Social Networking

Overview :

This assignment builds a social graph (tries to download the SNAP Facebook network, otherwise creates a synthetic SBM graph), computes per-node graph features, synthetically labels a small fraction of nodes as “bots”, trains a Random Forest classifier to detect bots, then simulates two adversarial attacks:

- **Structural evasion:** bots change their connectivity (remove some edges to high-degree neighbors and add edges to low-degree nodes) to try to evade detection.
- **Graph poisoning:** the training data is corrupted by flipping labels and injecting fake nodes with many edges; the model is retrained on this poisoned data and evaluated on the clean test set.



```
Assignment2 Social.py X graph_original.png
C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > graph_poisoning
1  import os
2  import gzip
3  import random
4  import urllib.request
5  import math
6  from collections import defaultdict
7  import argparse
8  import joblib
9
10 import networkx as nx
11 import numpy as np
12 import pandas as pd
13 import matplotlib.pyplot as plt
14
15 from sklearn.ensemble import RandomForestClassifier
16 from sklearn.model_selection import train_test_split
17 from sklearn.preprocessing import StandardScaler
18 from sklearn.pipeline import Pipeline
19 from sklearn.metrics import classification_report, accuracy_score, roc_auc_score
20
```

- Standard libs for files, randomness, downloading.
- networkx for graph operations.
- pandas/numpy for data handling.
- sklearn for modeling and metrics.
- joblib to save models.

```

16 from sklearn.model_selection import train_test_split
17 from sklearn.preprocessing import StandardScaler
18 from sklearn.pipeline import Pipeline
19 from sklearn.metrics import classification_report, accuracy_score, roc_auc_score
20
21 # -----
22 # Configurable parameters
23 # -----
24 SNAP_URL = "https://snap.stanford.edu/data/facebook_combined.txt.gz"
25 OUTPUT_DIR = "outputs"
26 RANDOM_SEED = 42
--

```



RANDOM_SEED used to make behavior reproducible (where supported). The script also sets Python and numpy seeds.

```

31
32 # -----
33 # Utilities
34 # -----
35 def ensure_outdir(path=OUTPUT_DIR):
36     os.makedirs(path, exist_ok=True)
37
38 def try_download_snap(save_to):
39     try:
40         print("[*] Attempting to download SNAP facebook_combined dataset...")
41         urllib.request.urlretrieve(SNAP_URL, save_to)
42         print("[+] Downloaded SNAP dataset to:", save_to)
43         return True
44     except Exception as e:
45         print("[!] Could not download SNAP dataset (internet blocked or error):", e)
46         return False
47
48 def load_snap_from_gz(path):
49     edges = []
50     with gzip.open(path, "rt") as f:
51         for line in f:
52             if not line.strip():
53                 continue
54             a,b = line.split()
55             edges.append((int(a), int(b)))
56     return edges
--

```



ensure_outdir(path=OUTPUT_DIR)

Creates the output directory if it doesn't exist.

try_download_snap(save_to)

Attempts to download the SNAP facebook dataset to save_to. Returns True on success, False on any exception (e.g., no internet).

load_snap_from_gz(path)

Reads the gzipped edge list line by line and returns a list of (a,b) integer edges. Used if the SNAP download succeeds.

```
54 |         a,b = line.split()
55 |         edges.append((int(a), int(b)))
56 |     return edges
57 |
58 | def build_sbm_graph(sizes=None, p_in=0.1, p_out=0.005):
59 |     """Build a synthetic social graph using stochastic block model (SBM)."""
60 |     if sizes is None:
61 |         sizes = [400,350,300,250,200,200,300] # ~2000 nodes
62 |         # Build probability matrix
63 |         nblocks = len(sizes)
64 |         p = [[p_in if i==j else p_out for j in range(nblocks)] for i in range(nblocks)]
65 |         G = nx.stochastic_block_model(sizes, p, seed=RANDOM_SEED)
66 |         G = nx.convert_node_labels_to_integers(G)
67 |         return G
68 |
```



build_sbm_graph(sizes=None, p_in=0.1, p_out=0.005)

- Builds a Stochastic Block Model (SBM) with sizes blocks (default block sizes sum to ~2000 nodes).
- Constructs a probability matrix p where intra-block edge prob = p_in and inter-block prob = p_out.
- Calls nx.stochastic_block_model and re-labels nodes to consecutive integers.
- Returns a networkx.Graph.

Use-case: fallback if SNAP download or load fails.

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > compute_graph_metrics

```

68
69 # -----
70 # Feature extraction
71 # -----
72 def compute_graph_metrics(G, approx_betweenness_k=200):
73     nodes = list(G.nodes())
74     print("[*] Computing degree and clustering...")
75     deg = dict(G.degree(nodes))
76     clustering = nx.clustering(G) # returns dict
77
78     # Betweenness: approximate when graph is large
79     n = G.number_of_nodes()
80     betweenness = {}
81     try:
82         if n > 1000:
83             # use k-approximation sampling for speed
84             k = min(approx_betweenness_k, n)
85             print(f"[*] Approximating betweenness centrality with k={k} samples...")
86             betweenness = nx.betweenness_centrality(G, k=k, seed=RANDOM_SEED)
87         else:
88             print("[*] Computing exact betweenness centrality...")
89             betweenness = nx.betweenness_centrality(G)
90     except Exception as e:
91         print("[!] Betweenness computation failed, setting zeros. Error:", e)
92         betweenness = {node: 0.0 for node in nodes}
93
94     # Eigenvector centrality using numpy
95     try:
96         print("[*] Computing eigenvector centrality (numpy)...")
97         eigen = nx.eigenvector_centrality_numpy(G)

```

Ln 73, Col 26 Spaces: 4 UTF-8 CRLF Python Signed out 3.13.9 (Microsoft Store)

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > compute_graph_metrics

```

72 def compute_graph_metrics(G, approx_betweenness_k=200):
96     print("[*] Computing eigenvector centrality (numpy)...")
97     eigen = nx.eigenvector_centrality_numpy(G)
98     except Exception as e:
99         print("[!] Eigenvector failed, falling back to degree-normalized proxy. Error:", e)
100         eigen = {node: float(deg.get(node, 0)) / (max(deg.values()) + 1e-6) for node in nodes}
101
102     print("[*] Computing average neighbor degree and core numbers...")
103     avg_neighbor_degree = nx.average_neighbor_degree(G)
104     try:
105         core_number = nx.core_number(G)
106     except Exception as e:
107         print("[!] core_number failed, filling zeros. Error:", e)
108         core_number = {node: 0 for node in nodes}
109
110     # Build DataFrame
111     data = {
112         'node': nodes,
113         'degree': [deg.get(n, 0) for n in nodes],
114         'clustering': [clustering.get(n, 0.0) for n in nodes],
115         'betweenness': [betweenness.get(n, 0.0) for n in nodes],
116         'eigenvector': [eigen.get(n, 0.0) for n in nodes],
117         'avg_neighbor_degree': [avg_neighbor_degree.get(n, 0.0) for n in nodes],
118         'k_core': [core_number.get(n, 0) for n in nodes],
119     }
120     df = pd.DataFrame(data).set_index('node')
121     return df
122
123 # -----
124 # Bot labeling (synthetic)

```

Ln 119, Col 6 Spaces: 4 UTF-8 CRLF Python Signed out 3.13.9 (Microsoft Store)

compute_graph_metrics(G, approx_betweenness_k=200)

For graph G, computes per-node features and returns a `pandas.DataFrame` indexed by node id. Features:

- `degree`: node degree (number of neighbors).
- `clustering`: clustering coefficient (triangles / possible triangles).
- `betweenness`: betweenness centrality. If $n > 1000$, uses k-sample approximate betweenness via `nx.betweenness_centrality(G, k=k, seed=RANDOM_SEED)` for speed; otherwise computes exact.
 - The code catches exceptions and falls back to zeros on failure.
- `eigenvector`: eigenvector centrality using `nx.eigenvector_centrality_numpy`. On failure falls back to a proxy: degree normalized by max degree.
- `avg_neighbor_degree`: average neighbor degree (from NetworkX).
- `k_core`: core number (`nx.core_number`). On failure fills zeros.

Important performance notes:

- Exact betweenness and eigenvector centralities can be expensive on large graphs; the script attempts approximate betweenness when $n > 1000$.

Return: `DataFrame` with one row per node and columns listed above.

```
121     return df
122
123     # -----
124     # Bot labeling (synthetic)
125     # -----
126     def label_synthetic_bots(df, fraction=0.05):
127         nodes = list(df.index)
128         n = len(nodes)
129         num_bots = max(1, int(fraction * n))
130         # pick half top-degree, half random
131         sorted_by_deg = df.sort_values('degree', ascending=False).index.tolist()
132         half = num_bots // 2
133         high_deg_bots = sorted_by_deg[:half]
134         remaining = list(set(nodes) - set(high_deg_bots))
135         random_bots = random.sample(remaining, num_bots - len(high_deg_bots))
136         bot_nodes = set(high_deg_bots + random_bots)
137         labels = pd.Series(0, index=df.index)
138         labels.loc[list(bot_nodes)] = 1
139         return bot_nodes, labels
140
```

label_synthetic_bots(df, fraction=0.05)

- fraction fraction of nodes are labeled as bots (default 5%).
- Strategy: choose half of bot nodes as the highest-degree nodes, and the other half randomly from the rest.
- Returns bot_nodes (a set) and labels (a pandas.Series indexed by nodes with 1 for bot, 0 otherwise).

This creates a plausible synthetic scenario where some bots are high-degree (influencer-like) and others look random.

```
139     return bot_nodes, labels
140
141     # -----
142     # Model training & eval
143     # -----
144     def train_evaluate(X_train, y_train, X_test, y_test, return_model=True):
145         clf = Pipeline([
146             ('scaler', StandardScaler()),
147             ('rf', RandomForestClassifier(n_estimators=200, random_state=RANDOM_SEED))
148         ])
149         clf.fit(X_train, y_train)
150         y_pred = clf.predict(X_test)
151         y_proba = clf.predict_proba(X_test)[:,1] if hasattr(clf, "predict_proba") else None
152         report = classification_report(y_test, y_pred, output_dict=True)
153         acc = accuracy_score(y_test, y_pred)
154         auc = roc_auc_score(y_test, y_proba) if y_proba is not None and len(set(y_test))>1 else float('na')
155         return {
156             'model': clf if return_model else None,
157             'report': report,
158             'accuracy': acc,
159             'auc': auc,
160             'y_pred': y_pred,
161             'y_proba': y_proba
162         }
```

train_evaluate(X_train, y_train, X_test, y_test, return_model=True)

- Builds a Pipeline with StandardScaler + RandomForestClassifier(n_estimators=200, random_state=RANDOM_SEED).
- Trains on training data, predicts on test data.
- Computes classification_report (as dict), accuracy, and AUC (if possible).
- Returns dictionary containing model, report, accuracy, auc, y_pred, y_proba.

Note: if all test labels are the same, ROC AUC is not computable and the code uses nan.

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > structural_evasion

```
162     }
163
164     # -----
165     # Attacks
166     # -----
167     def structural_evasion(G, bot_nodes, remove_frac=0.5, add_degree=5):
168         G2 = G.copy()
169         for b in bot_nodes:
170             if b not in G2:
171                 continue
172             nbrs = list(G2.neighbors(b))
173             if not nbrs:
174                 continue
175             # remove top-degree neighbors
176             nbrs_sorted = sorted(nbrs, key=lambda x: G2.degree[x], reverse=True)
177             num_remove = max(1, int(remove_frac * len(nbrs_sorted)))
178             to_remove = nbrs_sorted[:num_remove]
179             for r in to_remove:
180                 if G2.has_edge(b, r):
181                     G2.remove_edge(b, r)
182             # add edges to low-degree nodes
183             candidates = [n for n in G2.nodes() if n != b and not G2.has_edge(b, n)]
184             candidates_sorted = sorted(candidates, key=lambda x: G2.degree[x])
185             to_add = candidates_sorted[:add_degree] if len(candidates_sorted) >= add_degree else random.sample(candidates, add_degree)
186             for a in to_add:
187                 G2.add_edge(b, a)
188         return G2
189
```

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > graph_poisoning

```
167     def structural_evasion(G, bot_nodes, remove_frac=0.5, add_degree=5):
168         return G2
169
170     def graph_poisoning(G, train_nodes, flip_frac=0.15, inject_frac=0.02, edges_per_inject=20):
171         G2 = G.copy()
172         train_nodes = list(train_nodes)
173         num_flip = max(1, int(flip_frac * len(train_nodes)))
174         flipped = random.sample(train_nodes, num_flip)
175
176         # Inject nodes
177         n_inject = max(1, int(inject_frac * G.number_of_nodes()))
178         start_idx = max(G2.nodes() + 1 if len(G2.nodes()) > 0 else 0
179         injected = list(range(start_idx, start_idx + n_inject))
180         for node in injected:
181             G2.add_node(node)
182         # Connect injected nodes to many high-degree targets
183         high_deg = sorted(G2.nodes(), key=lambda x: G2.degree[x], reverse=True)[:200]
184         if not high_deg:
185             high_deg = list(G2.nodes())
186         for inj in injected:
187             targets = random.sample(high_deg, min(edges_per_inject, len(high_deg)))
188             for t in targets:
189                 G2.add_edge(inj, t)
190         return G2, flipped, injected
211
```

structural_evasion(G, bot_nodes, remove_frac=0.5, add_degree=5)

- For each bot b:
 - Get its neighbors and sort them by neighbor degree descending.
 - Remove the top `remove_frac` fraction of these neighbors (i.e., disconnect from high-degree neighbors).
 - Add `add_degree` new edges from b to low-degree nodes it wasn't connected to (candidates sorted by degree ascending). If there are not enough low-degree nodes, it picks a random subset.
- Returns a new graph G2 with the structural changes.
- Goal: make bot nodes less obviously high-degree/hub-linked and more blended with low-degree nodes.

graph_poisoning(G, train_nodes, flip_frac=0.15, inject_frac=0.02, edges_per_inject=20)

- Copies the graph to G2.
- Chooses `num_flip = flip_frac * len(train_nodes)` training nodes and marks them as flipped (this list is returned but not immediately applied to labels here — the caller flips labels in the training set).
- Injects `n_inject = inject_frac * number_of_nodes` new nodes with consecutive integer ids beyond current max node id.
- Connects each injected node to many high-degree targets (top 200 high-degree nodes) with up to `edges_per_inject` edges each.
- Returns (G2, flipped, injected).

Purpose: emulate poisoning the training graph by injecting noisy nodes and prepare to flip labels for some training nodes.

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > plot_graph_sample

```
211
212 # -----
213 # Visualization
214 # -----
215 def plot_graph_sample(G, bot_nodes=None, outpath=None, sample_size=500, title="Graph sample"):
216     if G.number_of_nodes() == 0:
217         return
218     if G.number_of_nodes() > sample_size:
219         sample = random.sample(list(G.nodes()), sample_size)
220         H = G.subgraph(sample).copy()
221     else:
222         H = G.copy()
223     pos = nx.spring_layout(H, seed=RANDOM_SEED)
224     plt.figure(figsize=(8,8))
225     if bot_nodes is None:
226         nx.draw_networkx_nodes(H, pos, node_size=20)
227     else:
228         node_color = []
229         for n in H.nodes():
230             node_color.append(1 if n in bot_nodes else 0)
231         nx.draw_networkx_nodes(H, pos, node_size=20, cmap=plt.cm.viridis, node_color=node_color)
232     nx.draw_networkx_edges(H, pos, alpha=0.3)
233     plt.title(title)
234     plt.axis('off')
235     if outpath:
236         plt.savefig(outpath, bbox_inches='tight', dpi=200)
237         plt.close()
238     else:
239         plt.show()
240
```

plot_graph_sample(G, bot_nodes=None, outpath=None, sample_size=500, title="Graph sample")

- If the graph has more nodes than sample_size, sample sample_size node IDs and work with the induced subgraph to keep drawing fast.
- Uses nx.spring_layout (seeded) to get positions and draws nodes and edges.
- When bot_nodes is given, nodes are colored by whether they're in the bot_nodes set (passed values 0/1 to node_color).
- Saves to outpath if provided, otherwise shows plot.

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > main_flow

```

242 # Main pipeline
243 # -----
244 def main_flow(force_synthetic=False):
245     ensure_outdir(OUTPUT_DIR)
246     snap_gz = os.path.join(OUTPUT_DIR, "facebook_combined.txt.gz")
247     G = None
248
249     if not force_synthetic:
250         downloaded = try_download_snap(snap_gz)
251         if downloaded:
252             try:
253                 edges = load_snap_from_gz(snap_gz)
254                 print("[*] Building graph from SNAP edges...")
255                 G = nx.Graph()
256                 G.add_edges_from(edges)
257             except Exception as e:
258                 print("[!] Failed loading SNAP edges, will build synthetic graph. Error:", e)
259                 G = None
260
261     if G is None:
262         print("[*] Building synthetic SBM graph (fallback)...")
263         G = build_sbm_graph()
264
265     print(f"[+] Graph: nodes={G.number_of_nodes()}, edges={G.number_of_edges()}")
266
267     # Compute baseline features
268     df_features = compute_graph_metrics(G)
269     df_features.to_csv(os.path.join(OUTPUT_DIR, "node_features_original.csv"))
270
271     # Label bots (synthetic)

```

Assignment2 Social.py X graph_original.png

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > main_flow

```

244 def main_flow(force_synthetic=False):
270
271     # Label bots (synthetic)
272     bot_nodes, labels = label_synthetic_bots(df_features, fraction=0.05)
273     df_features['is_bot'] = labels
274
275     # Split train/test (stratify)
276     FEATURE_COLS = ['degree', 'clustering', 'betweenness', 'eigenvector', 'avg_neighbor_degree', 'k_core']
277     X = df_features[FEATURE_COLS].fillna(0).values
278     y = df_features['is_bot'].values
279     node_ids = df_features.index.values
280
281     X_train, X_test, y_train, y_test, node_train, node_test = train_test_split(
282         X, y, node_ids, test_size=0.3, random_state=RANDOM_SEED, stratify=y
283     )
284
285     # Baseline model
286     print("[*] Training baseline classifier...")
287     baseline = train_evaluate(X_train, y_train, X_test, y_test, return_model=True)
288     joblib.dump(baseline['model'], os.path.join(OUTPUT_DIR, "baseline_model.joblib"))
289
290     # Save baseline report
291     with open(os.path.join(OUTPUT_DIR, "baseline_results.txt"), "w") as f:
292         f.write(classification_report(y_test, baseline['y_pred']))
293         f.write(f"\nAccuracy: {baseline['accuracy']}\nAUC: {baseline['auc']}\n")
294
295     print(f"[+] Baseline accuracy: {baseline['accuracy']}, AUC: {baseline['auc']}")
296
297     # Visualize original graph (sample)
298     plot_graph_sample(G, bot_nodes, outpath=os.path.join(OUTPUT_DIR, "graph_original.png"),

```

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > main_flow

```
244 def main_flow(force_synthetic=False):
288     joblib.dump(baseline['model'], os.path.join(OUTPUT_DIR, "baseline_model.joblib"))
289
290     # Save baseline report
291     with open(os.path.join(OUTPUT_DIR, "baseline_results.txt"), "w") as f:
292         f.write(classification_report(y_test, baseline['y_pred']))
293         f.write(f"\nAccuracy: {baseline['accuracy']}\nAUC: {baseline['auc']}\n")
294
295     print("[+] Baseline accuracy:", baseline['accuracy'], "AUC:", baseline['auc'])
296
297     # Visualize original graph (sample)
298     plot_graph_sample(G, bot_nodes=bot_nodes, outpath=os.path.join(OUTPUT_DIR, "graph_original.png"),
299
300     # -----
301     # Structural evasion
302     # -----
303     attacked_subset = random.sample(list(bot_nodes), max(1, int(0.3*len(bot_nodes)))) # attack 30% c
304     G_struct = structural_evasion(G, attacked_subset, remove_frac=0.5, add_degree=5)
305     df_struct = compute_graph_metrics(G_struct)
306     # Ensure ordering same as original nodes for comparison
307     df_struct = df_struct.reindex(df_features.index).fillna(0)
308     df_struct.to_csv(os.path.join(OUTPUT_DIR, "node_features_structural.csv"))
309
310     # Evaluate baseline model on features computed after structural evasion (no retraining)
311     X_struct_test = df_struct.loc[node_test, FEATURE_COLS].fillna(0).values
312     struct_eval = {}
313     # reuse baseline model to predict
314     model_baseline = baseline['model']
315     y_struct_pred = model_baseline.predict(X_struct_test)
316     try:
```

Ln 311, Col 26 Spaces: 4 UTF-8 CRLF {} Python Signed out 3.13.9 (Microsoft Store)

C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > main_flow > flipped_nodes

```
244 def main_flow(force_synthetic=False):
312     struct_eval = {}
313     # reuse baseline model to predict
314     model_baseline = baseline['model']
315     y_struct_pred = model_baseline.predict(X_struct_test)
316     try:
317         y_struct_proba = model_baseline.predict_proba(X_struct_test)[:,-1]
318     except Exception:
319         y_struct_proba = None
320     struct_report = classification_report(y_test, y_struct_pred, output_dict=False)
321     struct_acc = accuracy_score(y_test, y_struct_pred)
322     struct_auc = roc_auc_score(y_test, y_struct_proba) if y_struct_proba is not None and len(set(y_te
323
324     with open(os.path.join(OUTPUT_DIR, "structural_attack_results.txt"), "w") as f:
325         f.write(struct_report)
326         f.write(f"\nAccuracy: {struct_acc}\nAUC: {struct_auc}\n")
327
328     print("[+] After structural evasion -> accuracy:", struct_acc, "AUC:", struct_auc)
329     plot_graph_sample(G_struct, bot_nodes=bot_nodes, outpath=os.path.join(OUTPUT_DIR, "graph_structur
330
331     # -----
332     # Graph poisoning
333     # -----
334     # We'll poison training set: flip some labels and inject nodes, then retrain
335     G_poisoned, flipped_nodes, injected_nodes = graph_poisoning(G, node_train, flip_frac=0.15, inject
336     df_poison = compute_graph_metrics(G_poisoned)
337     df_poison = df_poison.reindex(df_features.index).fillna(0)
338     df_poison.to_csv(os.path.join(OUTPUT_DIR, "node_features_poison.csv"))
339
340     # Prepare poisoned training labels: flip labels for 'flipped_nodes' that are in train
```

```
Assignment2 Social.py X graph_original.png
C: > Users > shahd > OneDrive > Desktop > Assignment2 Social.py > main_flow
244 def main_flow(force_synthetic=False):
339
340     # Prepare poisoned training labels: flip labels for 'flipped_nodes' that are in train
341     # Find positions within node_train
342     y_train_poison = y_train.copy()
343     node_train_list = list(node_train)
344     for n in flipped_nodes:
345         if n in node_train_list:
346             idx = node_train_list.index(n)
347             y_train_poison[idx] = 1 - y_train_poison[idx] # flip
348
349     # Build poisoned training X using new features computed on poisoned graph
350     train_positions = [list(df_poison.index).index(n) for n in node_train] # mapping node id -> row
351     X_train_poison = df_poison.iloc[train_positions][FEATURE_COLS].fillna(0).values
352
353     # Retrain classifier on poisoned data
354     clf_poisoned = Pipeline([
355         ('scaler', StandardScaler()),
356         ('rf', RandomForestClassifier(n_estimators=200, random_state=RANDOM_SEED))
357     ])
358     print("[*] Training classifier on poisoned training data...")
359     clf_poisoned.fit(X_train_poison, y_train_poison)
360
361     # Evaluate poisoned-trained model on clean test features (from original df_features)
362     X_test_clean = df_features.loc[node_test, FEATURE_COLS].fillna(0).values
363     y_poisoned_pred = clf_poisoned.predict(X_test_clean)
364     try:
365         y_poisoned_proba = clf_poisoned.predict_proba(X_test_clean)[: ,1]
366     except Exception:
367         y_poisoned_proba = None
```

1. **Prepare outputs:** call `ensure_outdir(OUTPUT_DIR)`.
2. **Download / load graph:**
 - a. If `force_synthetic` is `False`, try to download the SNAP file to `outputs/facebook_combined.txt.gz`.
 - b. If download succeeds, load edges using `load_snap_from_gz` and create `G = nx.Graph()` with `G.add_edges_from(edges)`.
 - c. If any step fails, or if `force_synthetic=True`, fall back to `G = build_sbm_graph()` (the synthetic SBM).
3. **Print graph size:** nodes and edges.
4. **Compute baseline features:**
 - a. `df_features = compute_graph_metrics(G)` and save to `node_features_original.csv`.
5. **Label bots:**
 - a. `bot_nodes, labels = label_synthetic_bots(df_features, fraction=0.05)`.
 - b. Add labels as column `is_bot` in `df_features`.
6. **Prepare feature matrix and split:**

- a. FEATURE_COLS list:
`['degree', 'clustering', 'betweenness', 'eigenvector', 'avg_neighbor_degree', 'k_core']`.
 - b. `X = feature values, y = is_bot, node_ids = index (node ids)`.
 - c. `train_test_split(..., test_size=0.3, random_state=RANDOM_SEED, stratify=y)` returns `X_train, X_test, y_train, y_test, node_train, node_test`.
 - i. NOTE: stratify ensures similar bot fraction in both sets.
7. **Train baseline model:**
 - a. Calls `train_evaluate(...)` and saves model as `baseline_model.joblib`.
 - b. Writes a text file `baseline_results.txt` containing classification report, accuracy, and AUC.
 - c. Prints baseline accuracy and AUC.
8. **Visualize original graph:**
 - a. Calls `plot_graph_sample(..., outpath="outputs/graph_original.png")`.
9. **Structural evasion experiment:**
 - a. Select `attacked_subset` — 30% of `bot_nodes` at random to attack.
 - b. `G_struct = structural_evasion(G, attacked_subset, remove_frac=0.5, add_degree=5)`.
 - c. Compute features on `G_struct` → `df_struct` and reindex to original index.
 - d. Save `node_features_structural.csv`.
 - e. Evaluate the **baseline** model on the *evaded* features for the test nodes (no retraining).
 - i. Build `X_struct_test` from `df_struct.loc[node_test, FEATURE_COLS]`.
 - ii. `y_struct_pred = model_baseline.predict(X_struct_test)` and compute metrics.
 - f. Save `structural_attack_results.txt` and print the results.
 - g. Plot sample graph `structural_attack.png`.
10. **Graph poisoning experiment:**
 - a. `G_poisoned, flipped_nodes, injected_nodes = graph_poisoning(G, node_train, flip_frac=0.15, inject_frac=0.02, edges_per_inject=20)`.
 - b. Compute features on `G_poisoned` → `df_poison` and reindex.
 - c. Save `node_features_poison.csv`.
 - d. **Flip labels** in the training set:

- i. Copy `y_train` to `y_train_poison`.
 - ii. For every `n` in `flipped_nodes` that appears in `node_train`, flip the label at that index (`1 -> 0, 0 -> 1`).
 - e. Rebuild poisoned training `X_train_poison` using the new features `df_poison` for nodes in `node_train`.
 - i. The code maps nodes in `node_train` to their row positions in `df_poison` using `train_positions = [list(df_poison.index).index(n) for n in node_train]`.
 - f. Retrain a fresh RandomForest pipeline (same config) on (`X_train_poison`, `y_train_poison`).
 - g. Evaluate this poisoned-trained model on the **clean** test features `X_test_clean = df_features.loc[node_test, FEATURE_COLS]`.
 - h. Save `poisoning_attack_results.txt`, dump `poisoned_model.joblib`, and plot `graph_poisoning.png`.
 - i. Print results.
11. **Save a Markdown report:**
- a. `report_summary.md` contains counts and summarized metrics for baseline, structural evasion, and poisoning experiments and lists produced files.
12. **Finish:** prints that outputs are saved.

```
# -----
# CLI
# -----
def parse_args():
    p = argparse.ArgumentParser(description="Assignment 2: Bot detection + attacks pipeline")
    p.add_argument("--force-synthetic", action="store_true", help="Skip SNAP download and force synth")
    return p.parse_args()

if __name__ == "__main__":
    args = parse_args()
    main_flow(force_synthetic=args.force_synthetic)
```

`parse_args()` uses `argparse` to accept `--force-synthetic`.

```
rs\shahd\.vscode\extensions\ms-python.debugpy-2025.16.0-win32-x64\bundled\libs\debugpy\launcher' '59837' '--' 'c:\Users\shahd\OneDrive\Desktop\Assignment2 Social.py'
PS C:\Users\shahd\OneDrive\Desktop\outputs> & 'c:\Users\shahd\AppData\Local\Microsoft\WindowsApps\python3.13.exe' 'c:\Users\shahd\.vscode\extensions\ms-python.debugpy-2025.16.0-win32-x64\bundled\libs\debugpy\launcher' '59837' '--' 'c:\Users\shahd\OneDrive\Desktop\Assignment2 Social.py'
[*] Attempting to download SNAP facebook_combined dataset...
[+] Downloaded SNAP dataset to: outputs\facebook_combined.txt.gz
[*] Building graph from SNAP edges...
[+] Graph: nodes=4039, edges=88234
[*] Computing degree and clustering...
[*] Approximating betweenness centrality with k=200 samples...
[*] Computing eigenvector centrality (numpy)...
[*] Computing average neighbor degree and core numbers...
[*] Training baseline classifier...
[+] Baseline accuracy: 0.971947194719472 AUC: 0.7526475694444444
[*] Computing degree and clustering...
[*] Approximating betweenness centrality with k=200 samples...
[*] Computing eigenvector centrality (numpy)...
[*] Computing average neighbor degree and core numbers...
[+] After structural evasion -> accuracy: 0.9570957095709571 AUC: 0.7396122685185185
[*] Computing degree and clustering...
[*] Approximating betweenness centrality with k=200 samples...
[*] Computing eigenvector centrality (numpy)...
[*] Computing average neighbor degree and core numbers...
[*] Training classifier on poisoned training data...
[+] Poisoned training -> accuracy on clean test: 0.9645214521452146 AUC: 0.7565972222222221
[+] All outputs saved to: outputs
PS C:\Users\shahd\OneDrive\Desktop\outputs>
```

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	node	degree	clustering	betweenness	eigenvector	avg_neighbor	k_core								
2	0	347	0.041962	0.147377	3.31E-05	18.95965	21								
3	1	17	0.419118	1.56E-06	5.97E-07	48.23529	13								
4	2	10	0.888889	1.55E-07	2.17E-07	49.9	9								
5	3	17	0.632353	2.41E-06	6.56E-07	59.76471	13								
6	4	10	0.866667	0	2.17E-07	42.6	9								
7	5	13	0.333333	7.10E-06	1.18E-06	50.61538	10								
8	6	6	0.933333	0	2.11E-07	63.5	5								
9	7	20	0.431579	8.23E-05	2.57E-05	45.9	12								
10	8	8	0.678571	0	2.13E-07	48.375	5								
11	9	57	0.397243	6.81E-06	2.20E-06	42.40351	21								
12	10	10	0.822222	0	7.64E-07	79.1	10								
13	11	1	0	0	2.04E-07	347	1								
14	12	1	0	0	2.04E-07	347	1								
15	13	31	0.651613	8.21E-07	1.09E-06	54.54839	21								
16	14	15	0.742857	3.83E-07	2.24E-07	38.66667	10								
17	15	1	0	0	2.04E-07	347	1								
18	16	9	0.666667	0	2.55E-07	66.88889	9								
19	17	13	0.730769	3.54E-07	2.21E-07	42.76923	9								
20	18	1	0	0	2.04E-07	347	1								
21	19	16	0.283333	0	2.24E-07	33.125	7								
22	20	15	0.685714	3.61E-07	2.23E-07	37.4	9								
23	21	65	0.349038	0.001332	2.59E-05	42.4	21								
24	22	11	0.472727	0	1.03E-06	53.72727	9								

node_features_original

baseline_results.txt

FileEditView

	precision	recall	f1-score	support
0	0.97	1.00	0.99	1152
1	0.88	0.50	0.64	60
accuracy			0.97	1212
macro avg	0.93	0.75	0.81	1212
weighted avg	0.97	0.97	0.97	1212

Accuracy: 0.971947194719472

AUC: 0.7526475694444444

structural_attack_results.txt

File Edit View Close tab (Ctrl+W)

	precision	recall	f1-score	support
0	0.96	1.00	0.98	1152
1	1.00	0.13	0.24	60
accuracy			0.96	1212
macro avg	0.98	0.57	0.61	1212
weighted avg	0.96	0.96	0.94	1212

Accuracy: 0.9570957095709571

AUC: 0.7396122685185185

poisoning_attack_results.txt

×

+

File

Edit

View

	precision	recall	f1-score	support
0	0.97	0.99	0.98	1152
1	0.71	0.48	0.57	60
accuracy			0.96	1212
macro avg	0.84	0.74	0.78	1212
weighted avg	0.96	0.96	0.96	1212

Accuracy: 0.9645214521452146

AUC: 0.7565972222222221

Assignment2 Social.py

×

graph_original.png

×

□

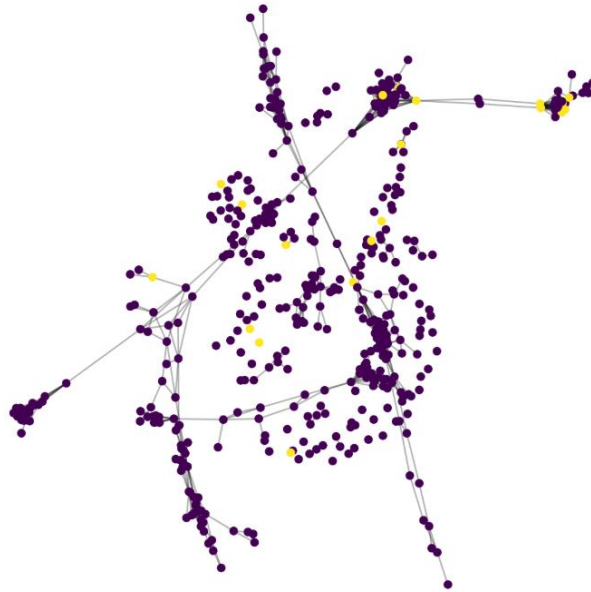
...

outputs > graph_original.png

Original Graph (sample)



After Structural Evasion (sample)



After Graph Poisoning (sample)

