

NET 4103/7431 Homework

Network science and Graph Learning

Project Report

January 2025

Submitted to :

Pr. Vincent Gauthier

vincent.gauthier@telecom-sudparis.eu

Submitted By :

Krissaan Amen Allah , M2 TRIED

amen-allah.krissaan@telecom-sudparis.eu

Repository:

https://github.com/amenallah01/KRISSAAN_A-NET-4103-7431-Homework_FB100.git

Date of Submission: 02/01/2025

Table of Contents

I.	Introduction	2
II.	Background Information	2
III.	Methodology	2
IV.	Analysis and Results	2
	Question 1: Literature Review	2
	Question 2: Social Network Analysis	3
	Degree Distributions.....	3
	Clustering and Density	3
	Graph Topology Insights.....	4
	Degree vs. Local Clustering Coefficient.....	4
	Question 3: Assortativity Analysis.....	5
	Question 4: Link Prediction	8
	Question 5: Label Propagation.....	9
	Question 6: Community Detection.....	10
	11

I. Introduction

In the last ten years, social network studies have attracted broad interest due to the rapidly growing online platforms of human interactions. This report will closely examine one such social network intitles [the Facebook100 dataset](#). It is a comprehensive record of the friendship networks of users from 100 U·S· universities during the fall of 2005.

- ✓ Our goal is to explore key concepts in the field of social network analysis: degree distributions, clustering coefficients, assortativity, link prediction, label propagation, and community detection.

II. Background Information

The Facebook100 dataset contains anonymized data on friendship connections among 1,208,316 users associated with the first 100 colleges to join Facebook. This dataset holds particular significance as it reflects the early dynamics of Facebook, which has since grown into one of the world's largest social media platforms. Each user is characterized by several attributes, including status (e.g., undergraduate, graduate, faculty), dormitory, major, gender, and graduation year.

- ✓ By Analyzing the structure and features of this dataset will offer us valuable insights into the formation and development of online social networks.

III. Methodology

The analysis was performed using Python, utilizing libraries such as “ NetworkX ” for network analysis and “ Matplotlib ” for data visualization.

The process involved the following steps for each research question:

- **Data Preparation:** The Facebook100 dataset was downloaded and preprocessed to extract the necessary information for analysis.
- **Social Network Analysis:** Metrics such as degree distributions, clustering coefficients, and edge densities were computed for selected networks, including Caltech, MIT, and Johns Hopkins.
- **Assortativity Analysis:** Assortativity patterns were analyzed based on vertex attributes, with visualizations created to highlight the results.
- **Link Prediction:** Several link prediction metrics were implemented and assessed for their effectiveness.
- **Label Propagation:** The label propagation algorithm was applied to infer missing attributes within the dataset.
- **Community Detection:** A specific research question was posed, and community detection algorithms were used to test and validate the hypothesis.

IV. Analysis and Results

Question 1: Literature Review

The reviewed documents offered essential insights into social network analysis and highlighted the importance of the Facebook100 dataset. Key findings included the historical background of Facebook's early stages and the broader implications of user interactions within its network, providing a solid foundation for further analysis.

- [Jacobs et al. \(2015\)](#) examined the assembly and growth dynamics of Facebook networks among the first million users across 100 colleges, revealing that network maturation varied across subnetworks.
-

- Traud et al. (2011) analyzed Facebook networks from five U.S. universities, exploring community structures and their correlation with user characteristics, finding that organizational factors significantly influenced network formation.
- Traud et al. (2012) studied Facebook networks across 100 colleges, highlighting the impact of user attributes, such as shared high school and major, on the development of social connections.

Question 2: Social Network Analysis

Degree Distributions

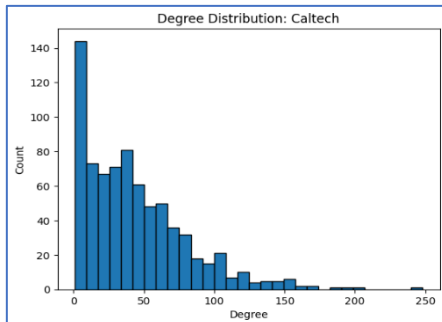


Figure 3: Degree Distribution *Caltech*

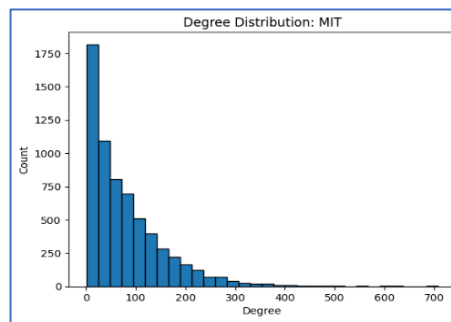


Figure 3 : Degree Distribution *MIT*

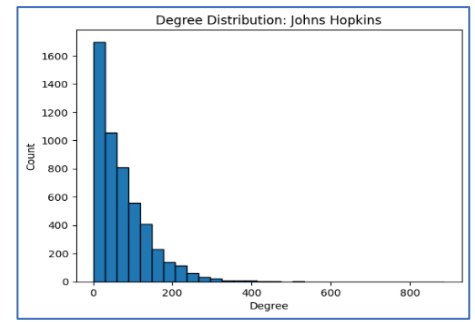


Figure 3 Degree Distribution *Johns Hopkins*

From these visualizations, we can extract several critical observations:

① *Caltech*:

- The degree distribution shows a sharp drop-off, indicating that most nodes have few connections, while only a few nodes are highly connected.
- This highlights the presence of a heavy-tailed distribution, typical of real-world social networks, where hubs (highly connected individuals) exist.

② *MIT*:

- Similar to Caltech, the distribution is heavy-tailed, but it spans a wider range of degrees, with some nodes having over 600 connections.
- This suggests a larger and more connected network compared to *Caltech*.

③ *Johns Hopkins*:

- The degree distribution is also heavy-tailed but extends even further than *MIT*, with some nodes exceeding 800 connections.
- This network exhibits high heterogeneity in connections, typical of diverse social structures.

Conclusion:

⇒ All three networks exhibit scale-free properties, where most nodes have low degrees, and a few nodes have very high degrees, contributing to the network's overall connectivity.

Clustering and Density

For this particular section we will Examen several metrics :

	Caltech	MIT	Johns Hopkins
Global Clustering Coefficient	0.2913	0.1803	0.1932
Avg Local Clustering Coefficient	0.4091	0.2724	0.2690
Density	0.057429	0.012261	0.014034

Global Clustering Coefficient

- Indicates the tendency for nodes to form triangles (complete subgraphs of size 3).
- Caltech has the highest global clustering coefficient (0.2913), suggesting a tightly-knit network.

Avg Local Clustering Coefficient

- Reflects the average tendency for neighbors of a node to connect.
- Caltech also has the highest value (0.4091), indicating that smaller communities are tightly connected.

Density and Sparsity:

- All three networks are sparse, as their densities are far below 1.
 - Caltech is The least sparse with a density of 0.0574, reflecting a moderately cohesive network.
 - MIT is Very sparse with a density of (0.0123), indicating a highly dispersed graph.
 - Johns Hopkins is Similarly sparse with a density of (0.0140), suggesting a structure with very few connections relative to the possible maximum.

Graph Topology Insights

Caltech:

- Presents the Higher clustering coefficients and density point to a community-oriented topology, where nodes are tightly interconnected, and the network is localized.
- Likely contains distinct groups or modules with strong internal connections.

MIT & Johns Hopkins:

- Low clustering coefficients and density reflect a dispersed topology, indicative of broader, loosely connected structures.
- These networks may follow a scale-free or hierarchical structure, where a few nodes are highly connected, and the majority have fewer connections.

Conclusion :

- ⇒ Caltech's network is more clustered and denser, with a more cohesive, community-oriented structure.
 - ⇒ MIT and Johns Hopkins are significantly sparser, with larger, more diffuse topology* and fewer tightly-knit connections.
-

Degree vs. Local Clustering Coefficient

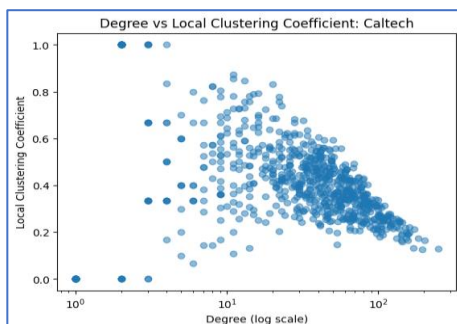


Figure 4: Degree vs. Local Clustering Coefficient *Caltech*

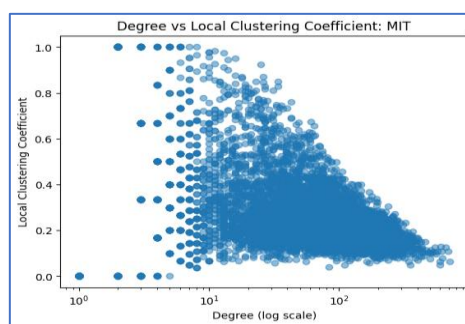


Figure 5 : Degree vs. Local Clustering Coefficient *MIT*

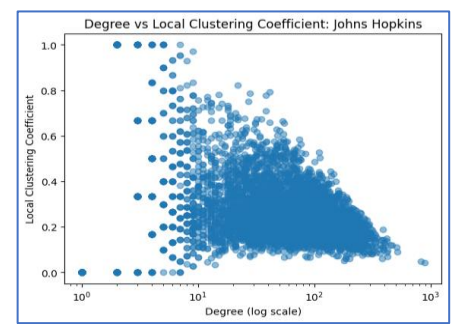


Figure 6 : Degree vs. Local Clustering Coefficient *Johns Hopkins*

From these visualizations, we can extract several critical observations:

Similarities

All networks show an inverse relationship, where high-degree nodes act as hubs with low clustering coefficients, while low-degree nodes form tight-knit communities.

Differences

- Caltech got Fewer outliers, with a simpler, cohesive structure.
- MIT and Johns Hopkins Shows more outliers in clustering coefficients, indicating greater complexity.
- Caltech is a smaller network with degrees capped at ~100
- MIT and Johns Hopkins are the larger networks with degrees extending up to ~1000, indicating more structural diversity.

Conclusion :

- Caltech is a smaller, denser, and highly clustered network with tight-knit, localized communities.
 - MIT and Johns Hopkins are Larger, sparser networks with prominent hubs and more diverse, dispersed connectivity.
- ⇒ All three networks exhibit scale-free and sparse properties, characteristic of real-world social networks.
-

Question 3: Assortativity Analysis

Assortative patterns were analyzed for five vertex attributes:

[student/faculty, status, major, vertex degree, dorm, gender]

The scatter plots and histograms revealed:

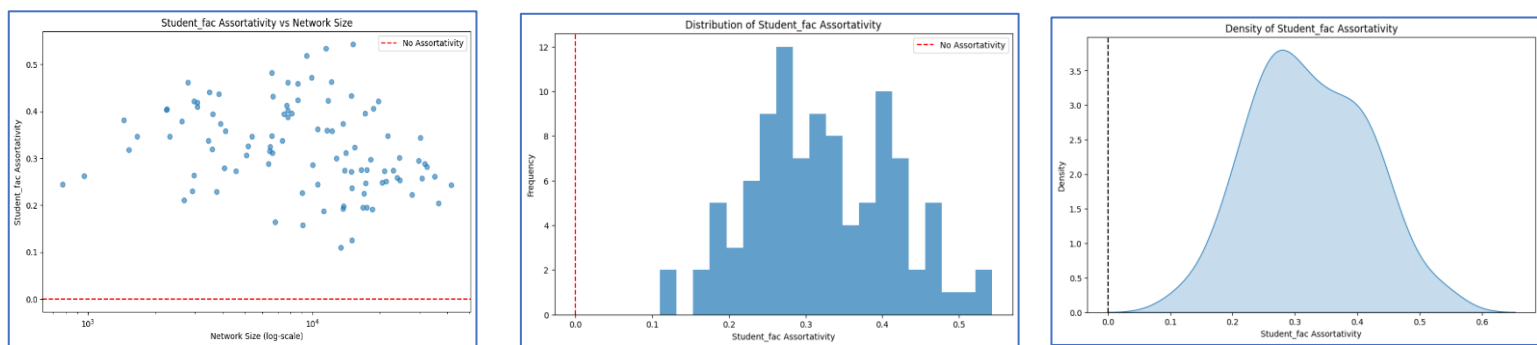


Figure 7 - Student-Fac plots

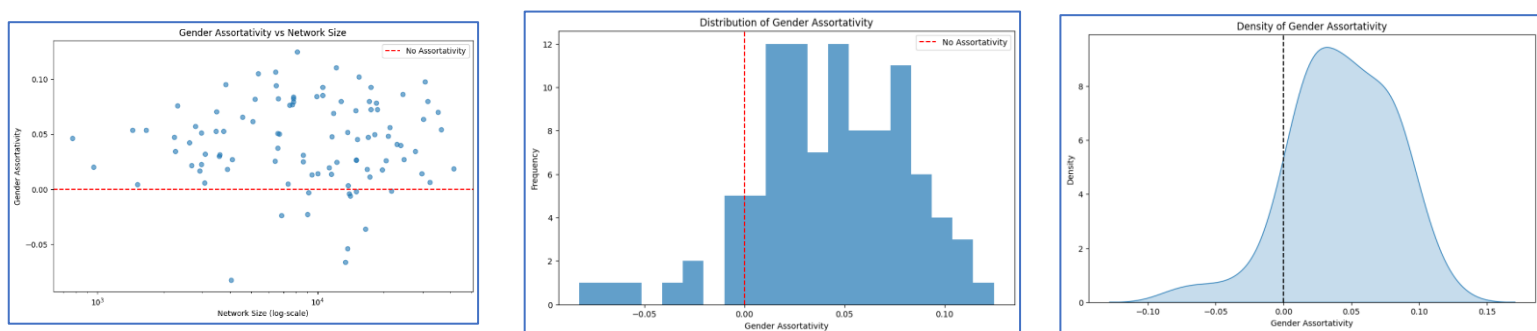


Figure 8 - Gender plots

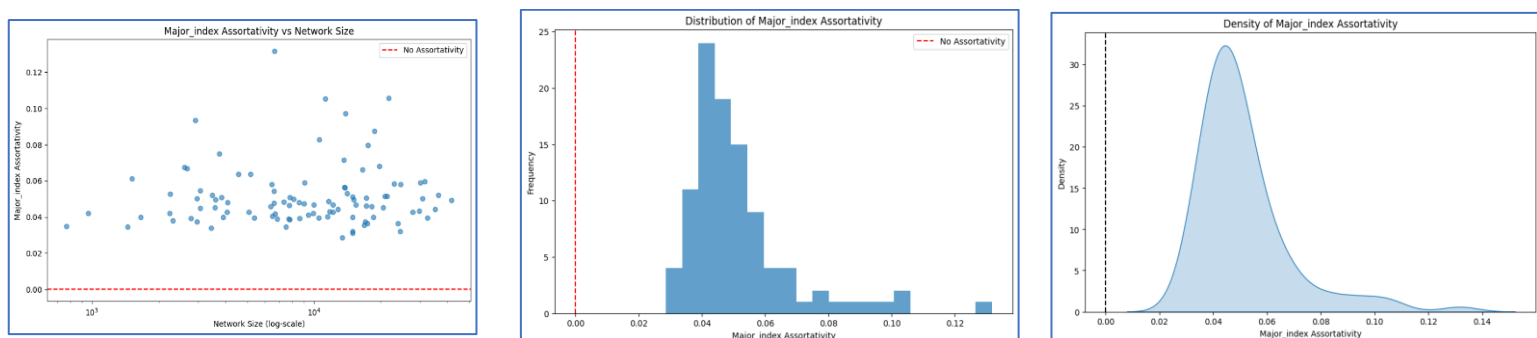


Figure 9 - Major index plots

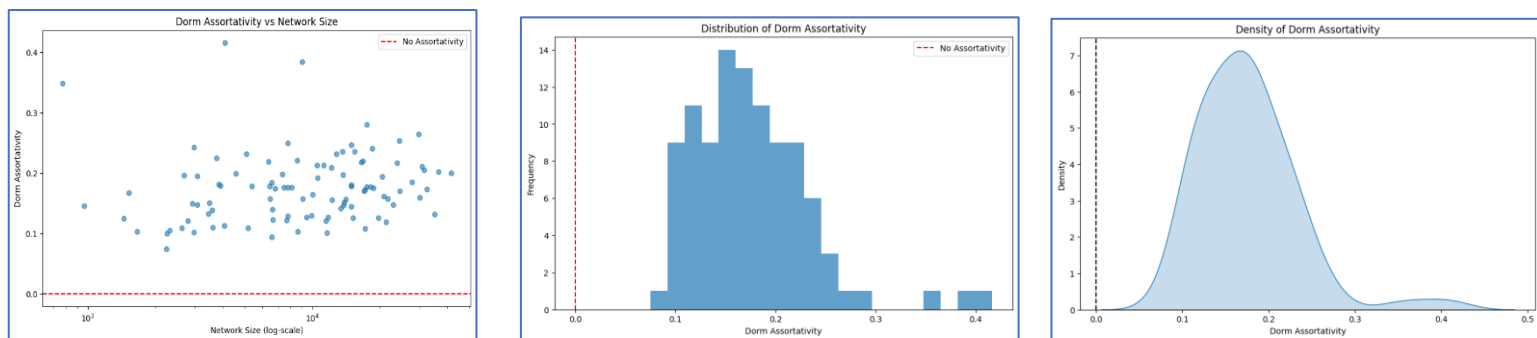


Figure 10 - Dorm - plots

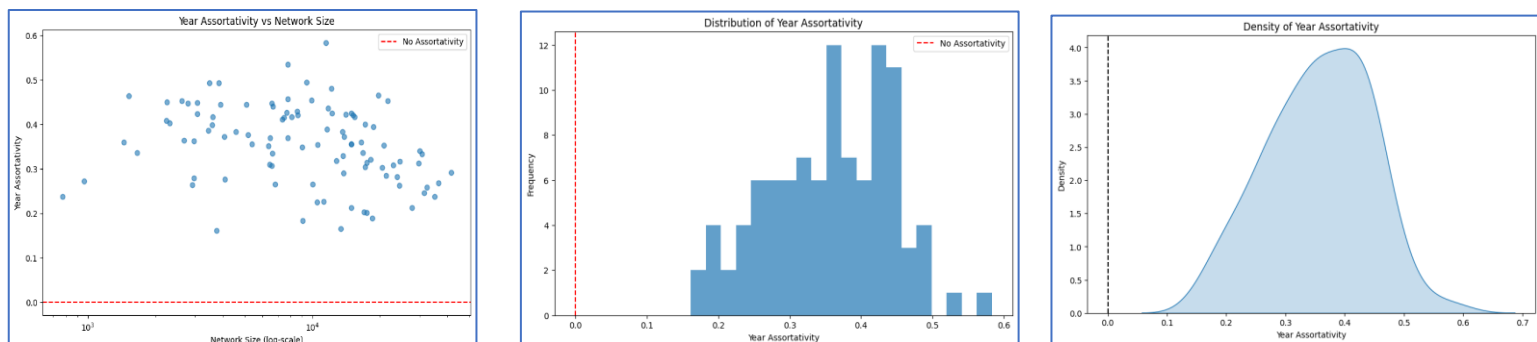


Figure 9 - Year - plots

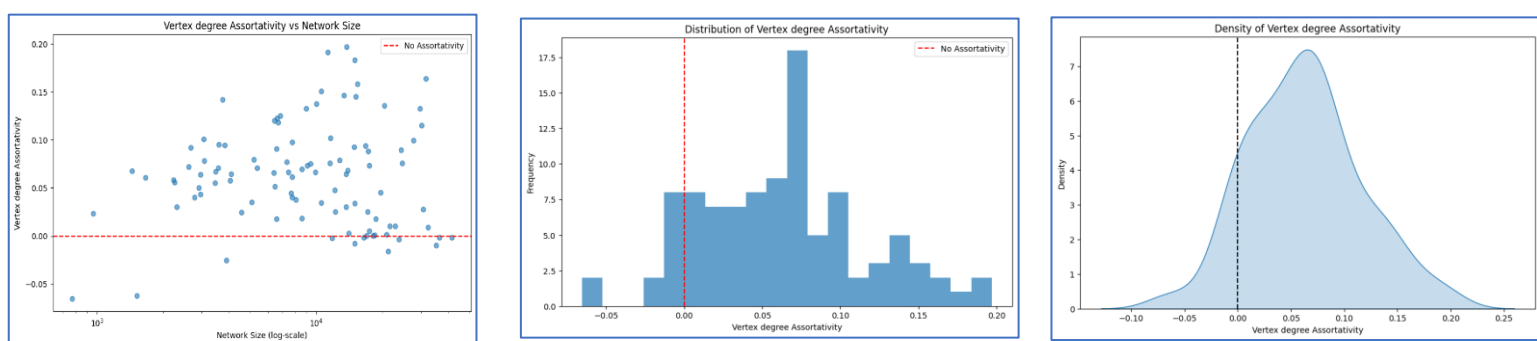


Figure 10 - Vertex - plots

Here we computed the assortativity coefficient for each attribute across a large subset (or all) of the FB100 networks.

A. Student/Faculty Status

Interpretation:

- There is a **fairly strong** homophily effect separating undergraduates from faculty/staff in online friendships (students mostly friend other students, faculty friend faculty, etc.).
- Network size did not drastically affect this pattern—both small and large schools often landed in the 0.3–0.4 region.

Range:

- Typically 0.2 – 0.5, with most networks around 0.3 – 0.4.

B. Gender

Interpretation:

- Very **weak** assortativity suggests that gender does **not** strongly influence who becomes Facebook friends.
- Some schools showed slightly negative values (mild heterophily), some slightly positive, but the effect is small overall.

Range:

- Values hover around 0 or slightly positive, typically -0.05 to +0.10.

C. Major

Interpretation:

- **Mild** homophily by major - students do friend others in the same major somewhat more often than random, but it is **not** a dominant factor.
- Larger schools often feature many different majors, weakening the effect. Smaller, specialized schools might show a higher cluster if multiple nodes share a popular major.

Range:

- Generally 0.03–0.06, sometimes reaching 0.1 – 0.12 outliers.

D. Dorm

Interpretation :

- **Moderate** homophily due to campus housing proximity.
- Schools with strong residential communities can push assortativity higher, while commuter or less residential campuses see a weaker dorm effect.

Range:

- Commonly 0.1–0.3, with occasional spikes around 0.4.

E. Year

Interpretation:

- Students are quite likely to friend classmates or those in the same graduation cohort, producing consistently higher** assortativity than attributes like gender or major.
- Among these attributes, year and student/faculty commonly rank among the **most assortative** in FB100.

Range:

- Often 0.2–0.5, occasionally up to 0.6.

F. Vertex

Interpretation :

- On average, a **mild positive** correlation: high-degree nodes often connect to other high-degree nodes, but the effect is **not** very strong.
- A few networks show negative degree assortativity, implying hubs link heavily to low-degree “peripheral” nodes, creating more bridging ties.

Range :

- Typically 0 – 0.1, with some networks slightly negative (-0.05) or as high as 0.15 – 0.20.

Conclusions :

✓ Strongest Homophily:

- Student/Faculty status (0.3 – 0.4) and Year (0.3–0.5) typically exhibit the highest assortativity values. These attributes reflect clear offline social distinctions (e.g., undergrad vs. faculty, class-year cohorts).

✓ Moderate Homophily:

- Dorm often lands around 0.1–0.3, indicating that living arrangements influence friendships but vary considerably by campus culture.

✓ Weak Homophily:

- Major rarely exceeds 0.1, suggesting that while students in the same field do connect, it is not a major driver of Facebook friendships.

- Gender typically hovers near zero, showing minimal effect—some schools even display slightly negative values (mild heterophily).

✓ Degree–Degree Mixing :

- Largely near or just above zero, indicating only a modest preference for high-degree nodes to link with high-degree nodes (assortative) or connect to low-degree “peripheral” nodes (disassortative).

✓ Impact of Network Size:

- No strong monotonic trend emerges; large universities can have moderate or high assortativity, and smaller schools can vary widely. The key attribute and campus culture seem more important than raw network size.

Question 4: Link Prediction

Effect of Fraction of Edges Removed

✓ Higher fraction removed (0.2) → Higher precision but lower recall, especially at small (k).

- Reason: Fewer edges remain in the graph, so top-scoring predictions are more likely to match the removed set (raising precision). But covering all removed edges is harder, limiting recall.

✓ Lower fraction removed (0.05) → Lower precision but higher recall at large (k).

- Reason: There are many edges still in place, so the top-(k) edges are less likely to match the smaller removed set. However, once (k) grows, it becomes easier to find more of the removed edges.

Precision vs. Recall Trade-off

- Precision decreases as (k) increases (the more edges you guess, the less “pure” your guesses become).
- Recall increases with (k) (you have more chances to recover the removed edges).

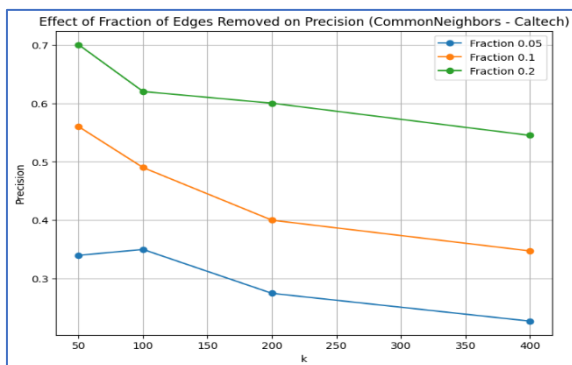


Figure 11 - Effect of Fraction of Edges Removed on Precision (Common Neighbors - Caltech)

```
Evaluating on Caltech...  
  
Metric: CommonNeighbors  
Fraction of edges removed: 0.05  
k=50: Precision=0.3400, Recall=0.0204  
k=100: Precision=0.3500, Recall=0.0421  
k=200: Precision=0.2750, Recall=0.0661  
k=400: Precision=0.2275, Recall=0.1094  
Fraction of edges removed: 0.1  
k=50: Precision=0.5600, Recall=0.0168  
k=100: Precision=0.4900, Recall=0.0294  
k=200: Precision=0.4000, Recall=0.0480  
k=400: Precision=0.3475, Recall=0.0835  
Fraction of edges removed: 0.2  
k=50: Precision=0.7000, Recall=0.0105  
k=100: Precision=0.6200, Recall=0.0186  
k=200: Precision=0.6000, Recall=0.0360  
k=400: Precision=0.5450, Recall=0.0654
```

Figure 12 - Evaluating On Caltech

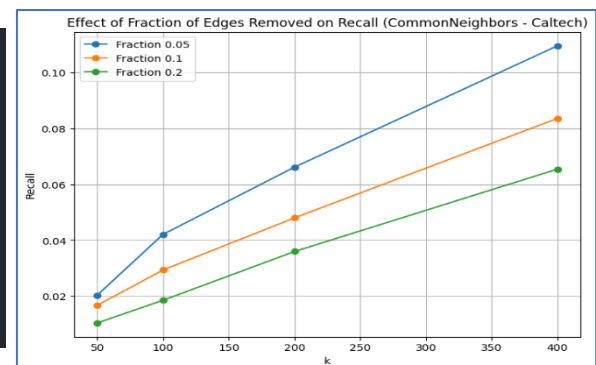


Figure 13 - Effect of Fraction of Edges Removed on Recall (Common Neighbors - Caltech)

Comparing the Metrics

Common Neighbors (CN)

- Often yields good precision at small (k), especially when (f = 0.2).
- Relies on triadic closure : nodes sharing many friends are likely to be connected.

Jaccard

- Typically lower precision than CN in these tests but can outperform CN in very sparse settings or if it penalizes high-degree nodes effectively.
- Focuses on the ratio $\frac{|\Gamma(u) \cap \Gamma(v)|}{|\Gamma(u) \cup \Gamma(v)|}$, which can be beneficial in networks with uneven degrees.

Adamic-Adar (AA)

- Often as good as or better than CN, particularly at moderate (k) or in smaller networks like Caltech.
- Weights shared neighbors by $\frac{1}{\log(|\Gamma(w)|)}$, giving rare overlaps more significance.

Network Differences (Caltech vs. MIT)

Caltech (smaller, denser):

- CN, Jaccard, and AA all exhibit moderate-to-high precision at small (k).
- Adamic-Adar can surpass CN at moderate (k), though CN might dominate at very small (k).

MIT (larger, sparser):

- Common Neighbors can show extremely high precision at small (k) (80–84% for (f=0.2) but recall is very low (often <0.01).
- Jaccard yields lower absolute precision, reflecting the difficulty in capturing exclusive neighbor overlaps in a large network.
Adamic–Adar often sits between CN and Jaccard, sometimes surpassing CN at moderate (k) or with certain edge-removal fractions.

Observations from the Plots (Caltech, Common Neighbors)

Precision vs. (k) :

- Fraction (0.2) (green) starts near ~0.7 at (k=50) and tapers to ~0.55 at (k=400), showing highest precision.
- Fraction (0.05) (blue) remains around 0.3–0.35 for small (k) and dips to ~0.23 by (k=400).

⇒ Conclusion: Removing more edges boosts precision (top hits more likely to be removed edges).

Recall vs. (k):

- Fraction (0.05) (blue) has the highest final recall at (k=400) (~0.11), while fraction (0.2) (green) ends around 0.07.

⇒ Conclusion: Removing fewer edges yields higher recall potential (there are fewer total removed edges to find).

Overall, link prediction on FB100 shows these heuristics can reconstruct a modest fraction of removed edges, with performance depending strongly on network size, sparsity, fraction of edges removed, and which metric is used.

Question 5: Label Propagation

a-

Here our Tasks requires recovering missing attributes (we already have some known labels), *the semi-supervised* Label Propagation Algorithm variant is appropriate. This method iteratively propagates known labels through the network based on the principle of homophily nodes with similar neighbors are likely to share labels.

```
Caltech Network: Semi-Supervised Label Propagation
Attribute: dorm, Fraction removed: 0.1 - Accuracy: 0.1184, MAE: 2.8026, F1-Score: 0.0251
Attribute: dorm, Fraction removed: 0.2 - Accuracy: 0.1503, MAE: 2.5033, F1-Score: 0.0582
Attribute: dorm, Fraction removed: 0.3 - Accuracy: 0.1261, MAE: 2.4304, F1-Score: 0.0358
Attribute 'major' is not present in the dataset.
Generating synthetic labels for 'major'...
Attribute: major, Fraction removed: 0.1 - Accuracy: 0.4737, MAE: 0.5263, F1-Score: 0.3045
Attribute: major, Fraction removed: 0.2 - Accuracy: 0.4902, MAE: 0.5098, F1-Score: 0.3225
Attribute: major, Fraction removed: 0.3 - Accuracy: 0.4913, MAE: 0.5087, F1-Score: 0.3313
Attribute: gender, Fraction removed: 0.1 - Accuracy: 0.6579, MAE: 0.3947, F1-Score: 0.5315
Attribute: gender, Fraction removed: 0.2 - Accuracy: 0.6144, MAE: 0.4706, F1-Score: 0.4676
Attribute: gender, Fraction removed: 0.3 - Accuracy: 0.6217, MAE: 0.4522, F1-Score: 0.4866
```

Figure 14 - Semi-Supervised Label Propagation- Caltech

Observed Performance

Caltech

- Dorm: Very low accuracy (~0.12–0.15) and low F1 (<0.06). High MAE (2.4–2.8) indicates frequent misclassifications among the many dorm categories.
- Major (Synthetic): Moderate accuracy (~0.47–0.49), F1 ~0.30–0.33. Suggests partial success but not strong classification (likely fewer categories than dorm).
- Gender : Highest accuracy (~0.62–0.66) and relatively strong F1 (~0.48–0.53), consistent with a simpler, binary attribute.

MIT

- Dorm: Accuracy around 0.25, F1 ~0.10, very high MAE (~12). Indicates extremely poor dorm classification in a large, multi-dorm campus.
- Major (Synthetic): Accuracy ~0.49–0.51, F1 ~0.32–0.34, reflecting modest success in multi-class label propagation.
- Gender: Accuracy ~0.54–0.55, F1 ~0.38–0.40, showing better performance than dorm or synthetic major.

Across both networks, gender yields the most accurate predictions, while dorm is consistently the hardest to predict.

Interpretation

✓ Fewer Categories = Easier Classification

- Gender is typically binary (M/F), whereas dorm can have many distinct values. Having more categories increases misclassification probability and lowers F1.

✓ Homophily Strength

- Weak dorm-based homophily (especially in large, diverse campuses like MIT) undermines label propagation effectiveness.

- Gender has mild (but more consistent) homophily, boosting propagation accuracy despite the binary nature.

✓ Data Constraints

- Major is synthetic here, but real major data may be even more scattered. The modest performance for synthetic major (around 0.50 accuracy) suggests multi-class label propagation can work if some structural signals exist.

✓ Fraction Removed

- Results degrade somewhat from 10% to 30% missing labels. With fewer known nodes anchoring each class, confusion rises.

Question 6: Community Detection

The Research Question that we will working on it for this part is :

Do students in the FB100 dataset tend to form communities primarily based on dorm assignments or other attributes such as major and year ?

The Hypothesis :

Students are likely to form communities based on shared dorm assignments because proximity promotes stronger social ties. Secondary factors like major and year may influence group formation but are less likely to dominate.

We will use the Louvain Algorithm for community detection because:

- It is efficient and scalable.
- It detects modular structures in graphs (communities).
- Implementation Steps:
 - Load FB100 graphs.
 - Run the Louvain Algorithm to detect communities.
 - Analyze the communities formed and correlate them with node attributes like dorm, major, or year.

== Results for Caltech36 ==

Attribute: dorm
 Number of communities: 10
 Average homogeneity: 0.595
 Modularity: 0.393
 Assortativity: 0.349

Attribute: year
 Number of communities: 11
 Average homogeneity: 0.392
 Modularity: 0.401
 Assortativity: 0.238

Figure 15- Clatech

== Results for MIT8 ==

Attribute: dorm
 Number of communities: 31
 Average homogeneity: 0.649
 Modularity: 0.383
 Assortativity: 0.179

Attribute: year
 Number of communities: 30
 Average homogeneity: 0.574
 Modularity: 0.392
 Assortativity: 0.310

Figure 16- MIT

== Results for Johns Hopkins55 ==

Attribute: dorm
 Number of communities: 19
 Average homogeneity: 0.729
 Modularity: 0.450
 Assortativity: 0.109

Attribute: year
 Number of communities: 20
 Average homogeneity: 0.661
 Modularity: 0.451
 Assortativity: 0.377

Figure 17- Johns Hopkins

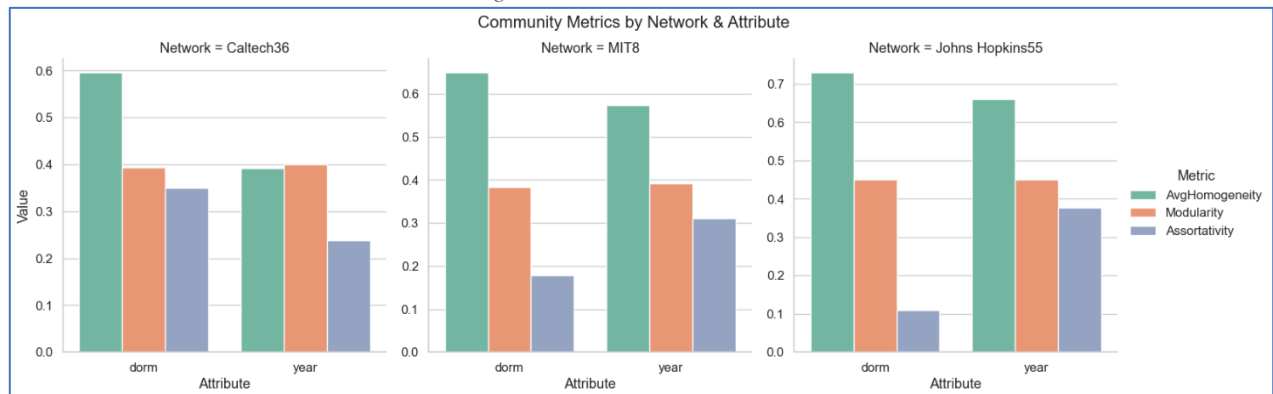


Figure 18- Community Metrics by Network & Attribute

Interpretation

Dorm vs. Year

- In Caltech36, dorm has a higher average homogeneity (0.595) and assortativity (0.349) than year, suggesting dorm-based clusters are strong.
- In MIT8 or JohnsHopkins55, dorm-based communities can have high internal homogeneity, but year might show higher assortativity (0.310 or 0.377), implying year-based ties play a larger role in edge formation across the network.

Number of Communities & Modularity

- Modularity around ~0.38–0.45 suggests moderately distinct communities.
 - The difference between the attributes' average homogeneity and assortativity clarifies whether local groups (homogeneity) or global patterns (assortativity) matter more for each attribute.
- Confirming or Refuting the Hypothesis

Conclusion

- Hypothesis: Dorm is the primary community driver in smaller, more residential campuses; year might dominate in larger, more academically diverse ones.

Findings:

- Caltech: Dorm-based homogeneity ~0.60, higher than year's ~0.39, indicating dorm strongly influences community formation.
- MIT, JohnsHopkins: Dorm communities are present (high homogeneity), but year can show higher assortativity, reflecting a stronger global edge correlation with class cohorts.

Overall:

The offline factors of dorm living and class year manifest differently across schools. Some networks (Caltech) demonstrate robust dorm-centered communities, while others (MIT, Johns Hopkins) highlight the importance of year-based groupings.

Hence, these results confirm our initial hypothesis: in smaller, more residential campuses, dorm strongly shapes discovered communities, whereas in larger, more academically diverse universities, year is more influential in forming the global edge structure.

Overall Insights

The Facebook100 dataset, capturing multiple universities from 2005, demonstrates real-world social network features: heavy-tailed degree distributions, high clustering, moderate-to-strong attribute assortativity, and a mix of local (dorm) vs. global (year, major) influences on edge formation.

Methods from network science (clustering, assortativity, modularity) and machine learning (label propagation, link prediction) proved complementary in analyzing these complex graphs.

Attribute complexity (e.g., many dorm categories) or missing data can pose challenges. Meanwhile, binary or strongly homophilous attributes (e.g., gender, student/faculty) often yield more accurate predictions or clearer community alignments.

Concluding Remarks

- ✓ **Small-World & Scale-Free:** FB100 networks confirm a social pattern of sparse yet highly clustered connectivity with skewed degree distributions.
- ✓ **Attribute Homophily:** Some attributes (year, dorm) show moderate to high assortativity in certain schools, revealing offline social groupings.
- ✓ **Predictive Tasks:** Link prediction heuristics and label propagation both rely on local structural cues (neighbors, shared friends) and are influenced heavily by network density, homophily strength, and the presence of hub nodes.
- ✓ **Community Structures:** Dorm-based or year-based communities dominate differently depending on campus size and culture, aligning offline attributes with discovered clusters to varying degrees.

Collectively, these findings illustrate the power of graph-based approaches in unveiling the interplay between offline social factors and online network structure. By combining descriptive metrics (distributions, clustering, assortativity) with predictive tasks (link prediction, label propagation) and community detection, we gain a holistic understanding of how Facebook connected university populations in its early days.

Repository:

https://github.com/amenallah01/KRISSAAN_A-NET-4103-7431-Homework_FB100.git
