# How community-like is the structure of synthetically generated graphs?

**Arnau Prat-Pérez**
**Universitat Politècnica de Catalunya**
**Barcelona**

David Dominguez-Sal
Sparsity Technologies
Barcelona

DAMA-UPC. DATA MANAGEMENT
UNIVERSITAT POLITÈCNICA DE CATALUNYA

*Sparsity

LDBC

# Motivation

- Community Detection is typically tested using **synthetic graphs (LFR generator)**.

  – Not only the graph output, but **communities** also.

- Recently, real graphs with ground truth have acquired popularity.

- How realistic is the community structure of synthetically generated graphs?

  – Existing work on vertex centric characteristics.

# Methodology

- We select real datasets with **ground truth communities**.
- We select two synthetic generators: **LFR** and **LDBC Data Generator**.
  - They output communities.
- We select a set of **6 metrics**.
- For each pair of graphs and each metric, we compare the distributions of the communities using the **Spearman's correlation coefficient**.

# Real Graphs

- Widely used in the literature.
- Diverse origin.
- Different sizes.

|  | Nodes | Edges |
|---|---|---|
| Amazon | 334,863 | 925,872 |
| Dblp | 317,080 | 1,049,866 |
| Youtube | 1,134,890 | 2,987,624 |
| LiveJournal | 3,997,962 | 34,681,189 |

# LFR Generator

- LFR

    - Generator created as a benchmark for Community Detection.
    - Five graphs with different mixing factors: 0.1 to 0.5.
    - Other parameters matching those found in real graphs.
    - Communities directly output by the program.

| | Nodes | Edges |
|---|---|---|
| Lfr.1 | 150,000 | 649,538 |
| Lfr.2 | 150,000 | 650,163 |
| Lfr.3 | 150,000 | 650,946 |
| Lfr.4 | 150,000 | 649,363 |
| Lfr.5 | 150,000 | 648,128 |

# LDBC Data Generator

- LDBC Data Generator
  - Data Generator of the LDBC Social Network Benchmark.
  - Communities are created from metadata.
  - One instance, simulating 3 years of 150000 users activity.

|  | Nodes | Edges | Communities |
|---|---|---|---|
| LDBC | 150,000 | 5,530,880 | 2,110,508 |

# Metrics

- 4 metrics for the internal structure:

  - Clustering Coefficient

  - Triangle Participation Ratio (TPR)

  - Bridge Ratio

  - Diameter

- 1 metric for the external connectivity.

  - Conductance

- Also the Size.
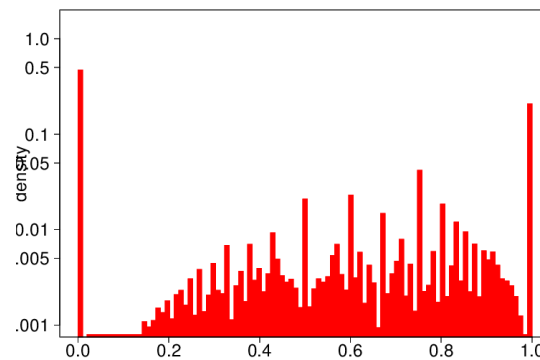
# Correlations



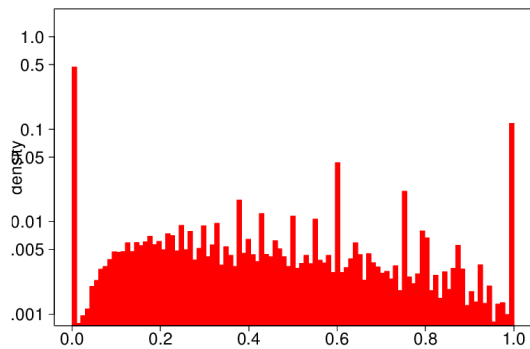Clustering Coefficient



TPR



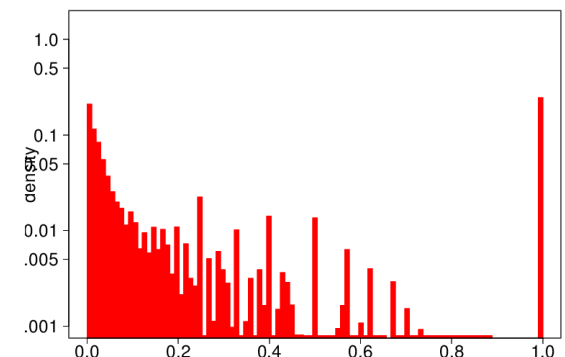Bridges Ratio
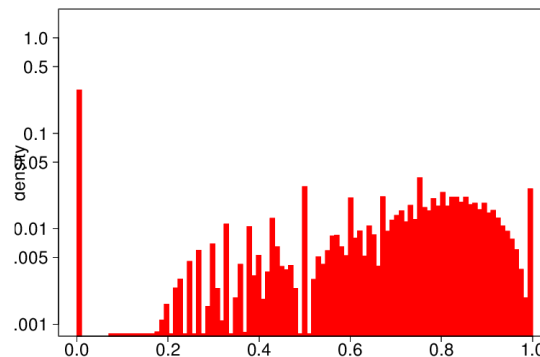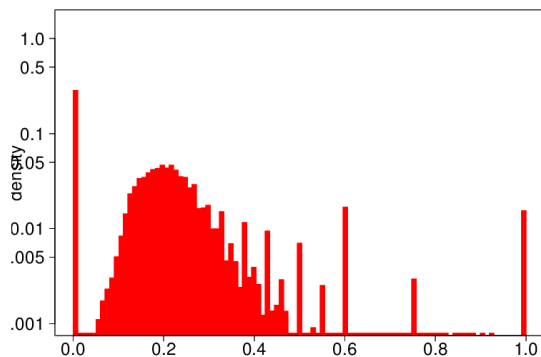


Log10(Size)

8

# Multimodality

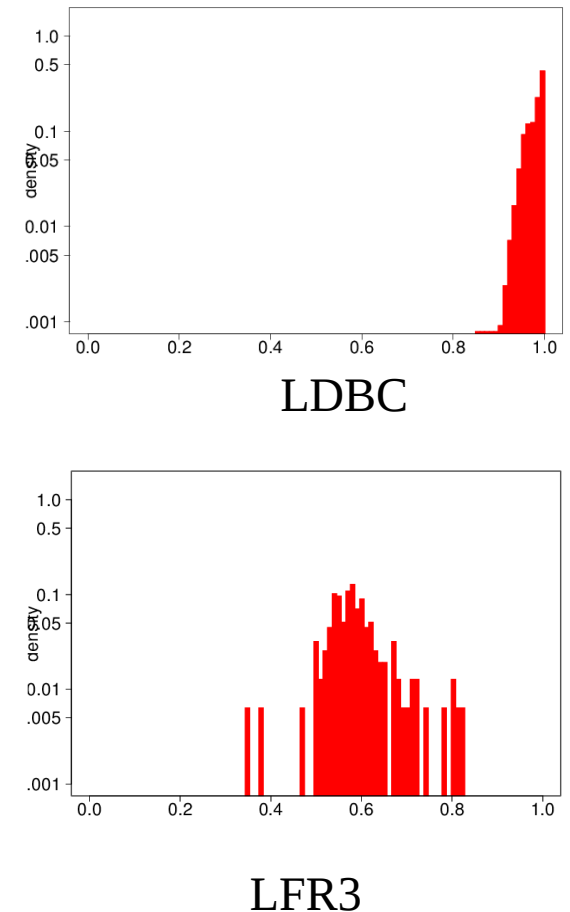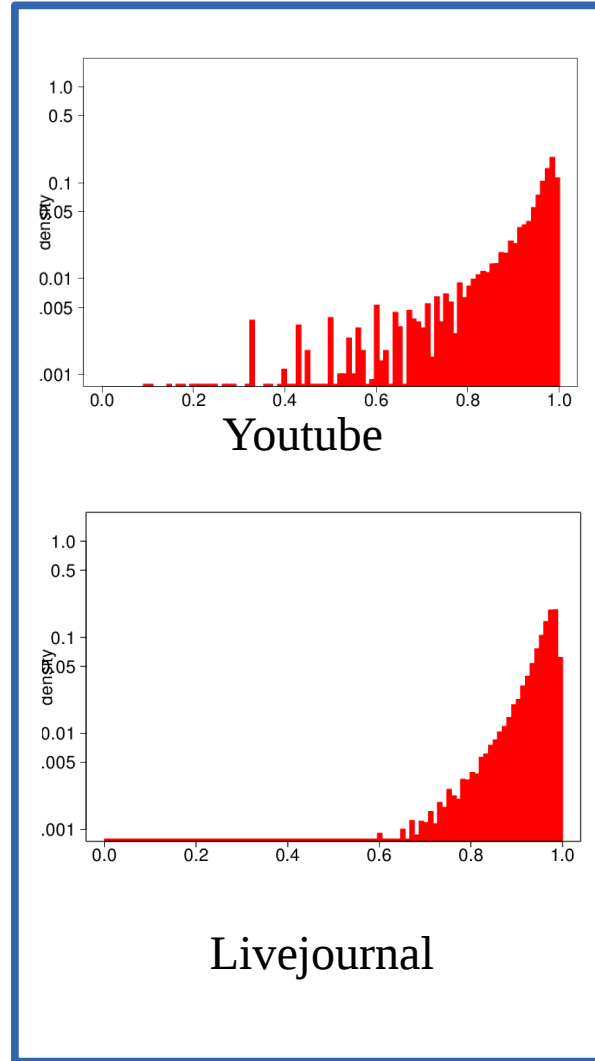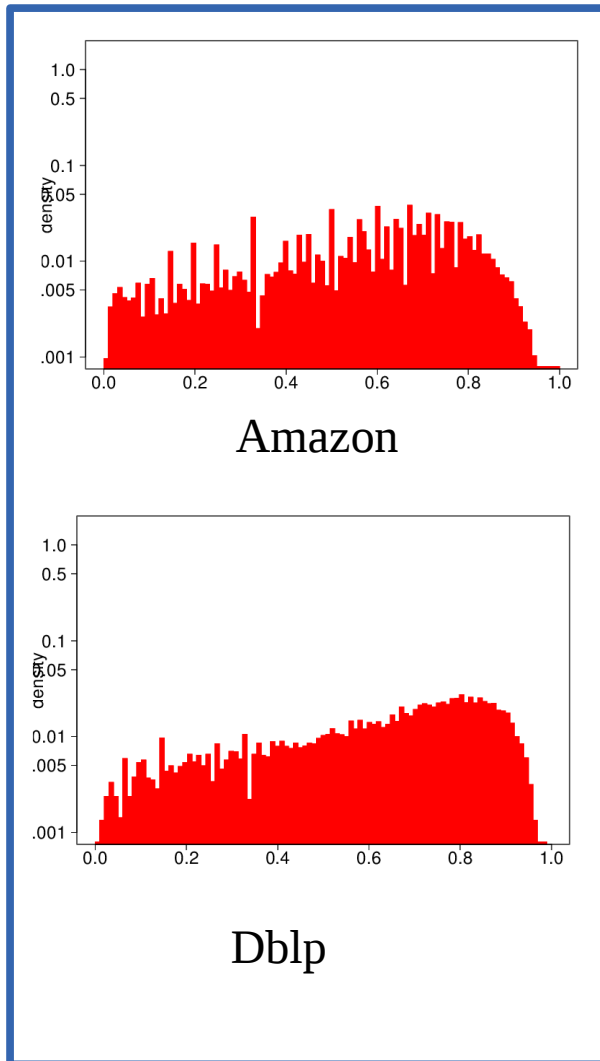- Multimodal distributions for CC, TPR and Bridge Ratio.



LiveJournal

LDBC

Clustering Coefficient          TPR          Bridge Ratio

# Findings on real graphs

- Signs of two different **Conductance** profiles



Amazon

Youtube

LDBC

Dblp

Livejournal

LFR3

# Conclusions

- Real graphs show similar distributions.

- LDBC Data Generator distributions are more realistic than those produced by LFR.

- Some distributions are multimodal: LDBC Data Generator mimics this.

- Signs of two different conductance profiles.

- Future Work: Experiment with more parameter configurations.

Thank you!