

463.2 Social Networks

CS463/ECE424

University of Illinois



Homophily in Social Networks
Social Network Inference
Privacy Risks
Discussion





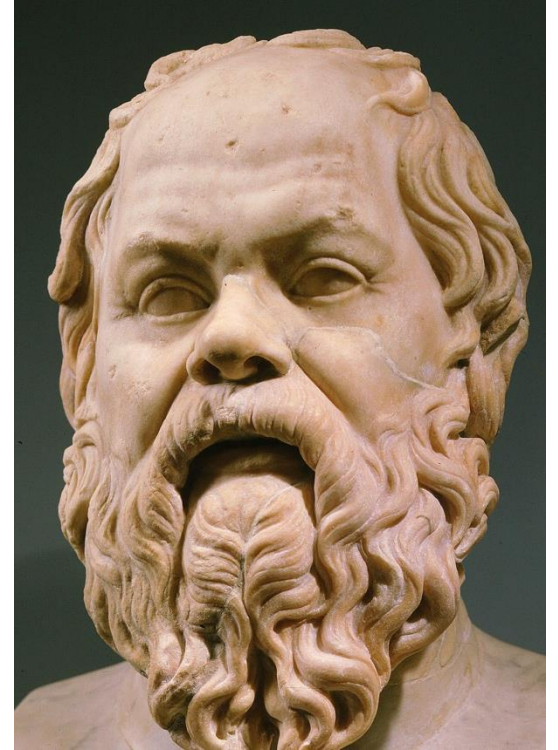
People of similar characteristics tend to befriend each other



463.2.1 Homophily in Social Networks

Homophily

- **Homophily** (i.e., "love of the same") is the tendency of individuals to associate and bond with similar others.
 - Term coined in 1950s in sociology papers.
- Systematically studied even earlier
- Much older concept; **Socrates to Lysis**:
 - “... and have you not also met with the treatises of philosophers who say that like must love like ...”
- **Modern variant**: ‘Similarity breeds connection’



Homophily

- Shown to exist for many attributes
 - Race/Ethnicity
 - Age
 - Religion
 - Education
 - Occupation
 - Gender
 - Marriage (homogamy)

Socrates speaking to a pair of youths:

I shall not ask which is the richer of the two, I said; for you are friends, are you not?

Certainly, they replied.

And friends have all things in common, so that one of you can be no richer than the other, if you say truly that you are friends. They assented.

Homophily: Terminology

Choice Homophily

Closeness due to preferences by the individual.

Example: Favorite teams

Induced Homophily

Closeness due to other constraints.

Examples: Geographic closeness, Age closeness with friends.

Value Homophily

Individuals with similar values, thinking.

Example: Religion

Status Homophily

Individual with similar social status.

Example: Aristocracy

Geographic Homophily: Marriages

- George Zipf, studied a large number of such empirical relationships

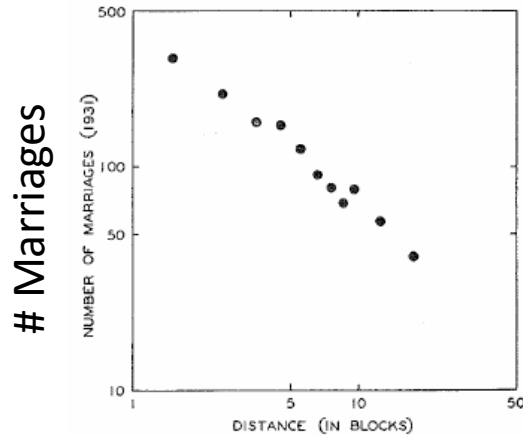


Fig. 9-22. Number of marriage licenses issued to 5,000 pairs of applicants living within Philadelphia in 1931 and separated by varying distances (the data of J. H. S. Bossard).

- Is this (inverse) relationship independent of other factors?

Geographic Homophily

- Size and distance of populations correlate with their degree of connection.

- Zipf equation:
Connection =

$$G * ((Pop1 * Pop2) / Distance)$$

G is a scaling factor

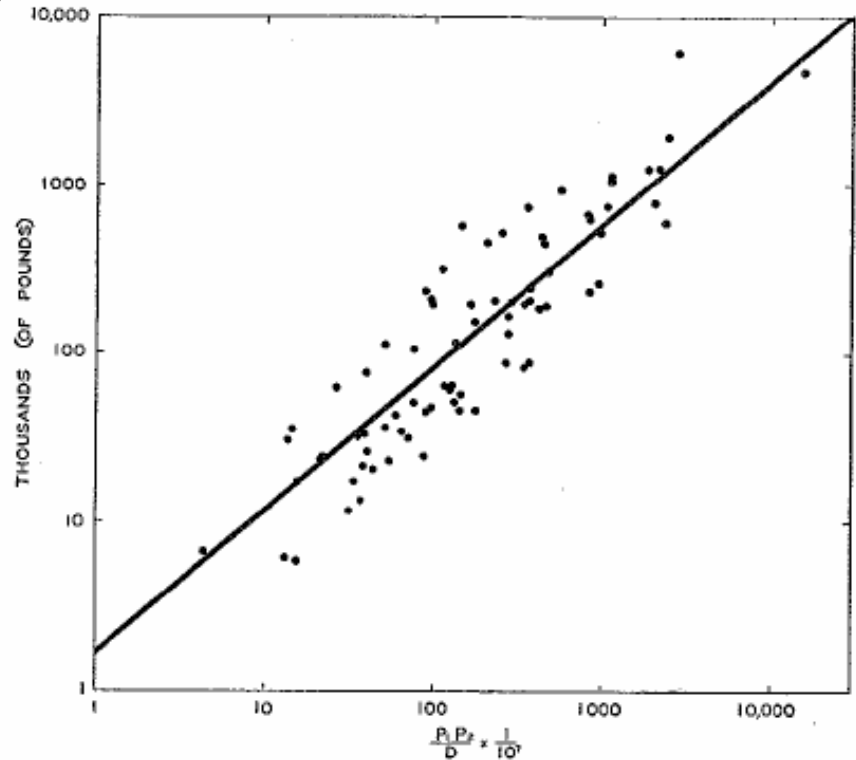
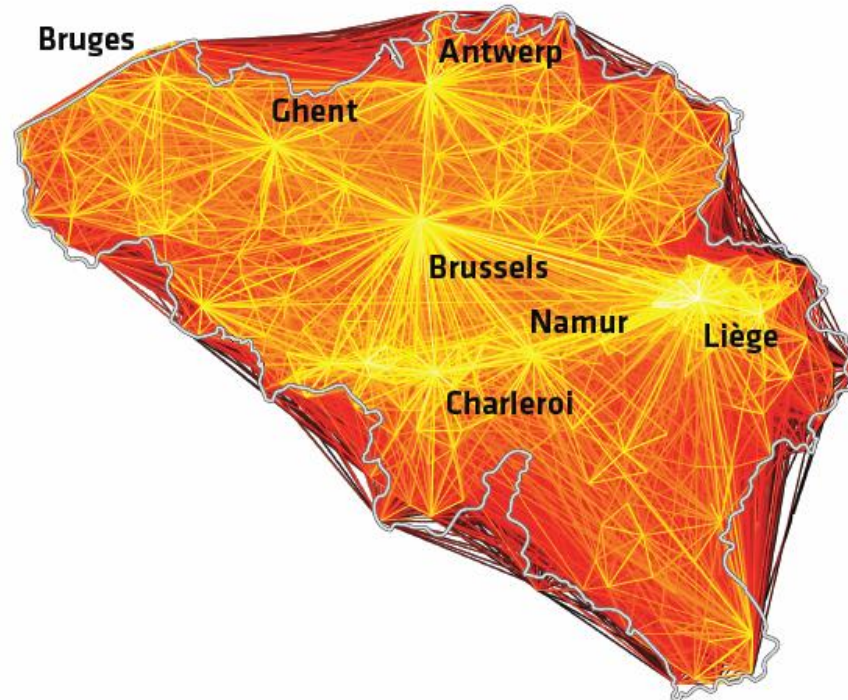


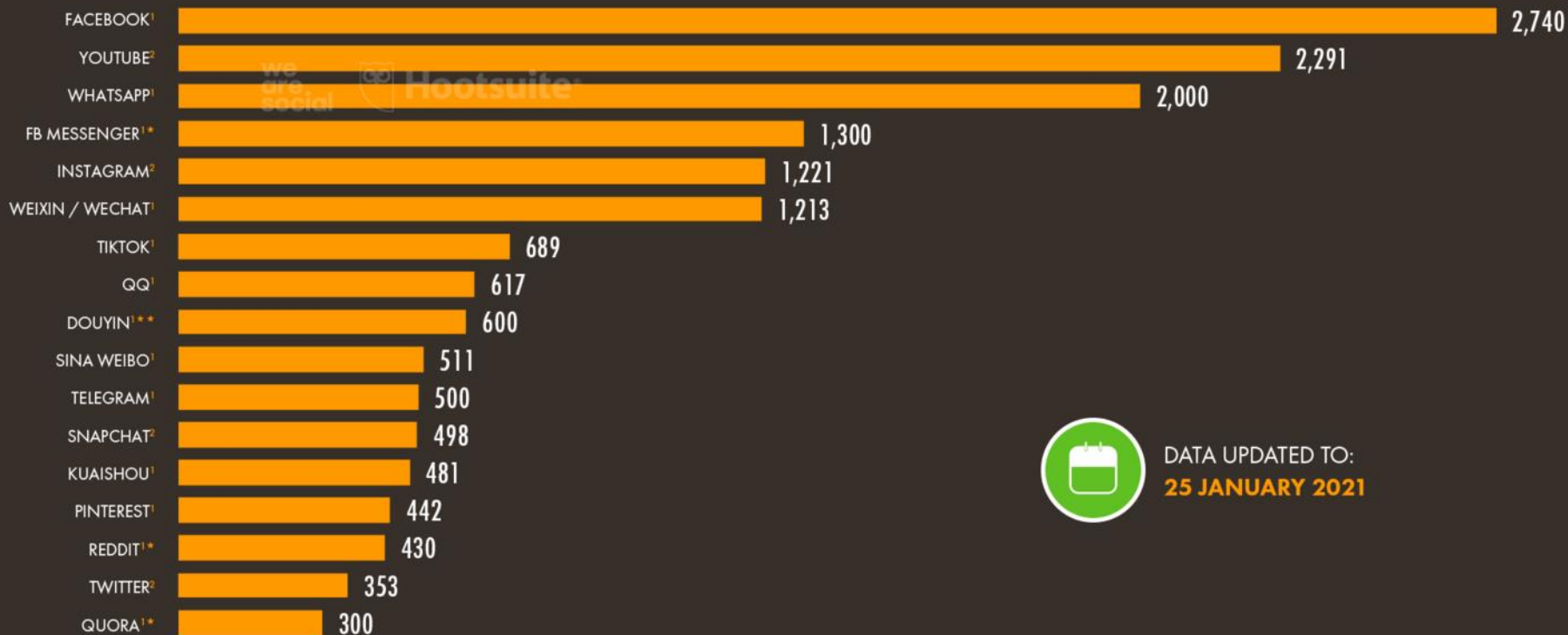
Fig. 9-14. Railway express. The movement by weight (less carload lots) between 13 arbitrary cities in the U. S. A., May 1939.

Geographic Homophily: Telephone Call Graphs in Belgium



THE WORLD'S MOST-USED SOCIAL PLATFORMS

THE LATEST GLOBAL ACTIVE USER FIGURES (IN MILLIONS) FOR A SELECTION OF THE WORLD'S TOP SOCIAL MEDIA PLATFORMS*



DATA UPDATED TO:
25 JANUARY 2021

Milgram's Six Degrees of Separation (Small-World)

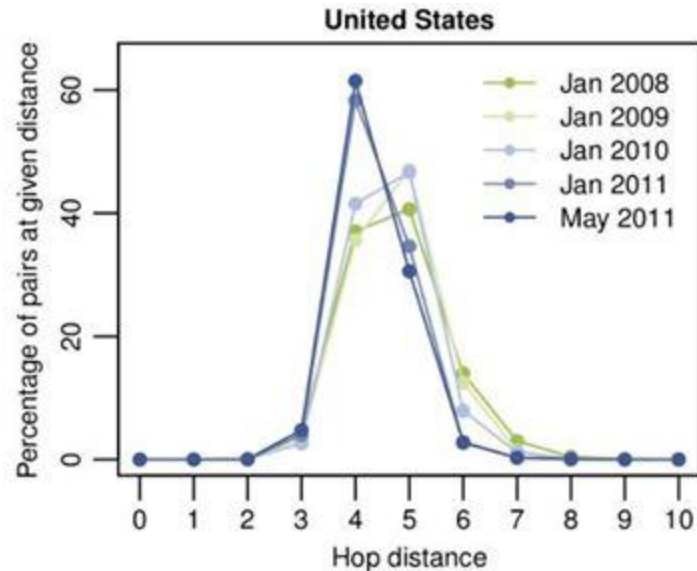
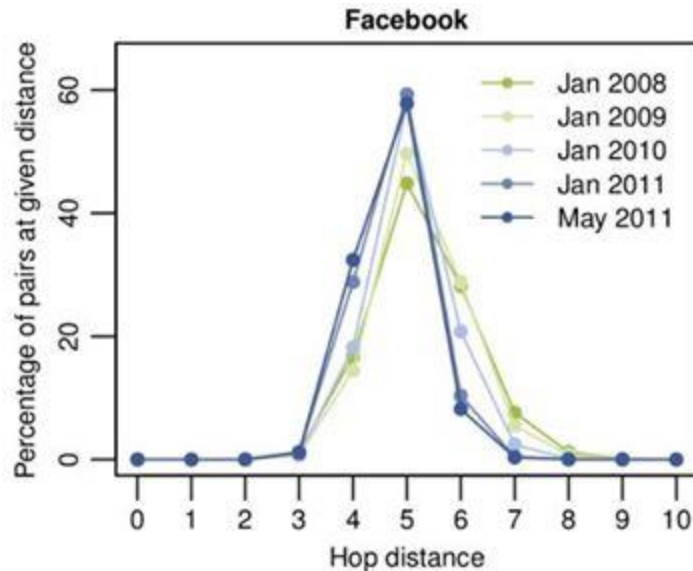
- The Six Degrees of Separation (Milgram 1967)
- Random people from Nebraska were to send a letter (via intermediaries) to a stockbroker in Boston.
- Could only send to someone with whom they were on a first-name basis.
- Not many arrived, but among the letters that found the target, the average number of links was **six**.



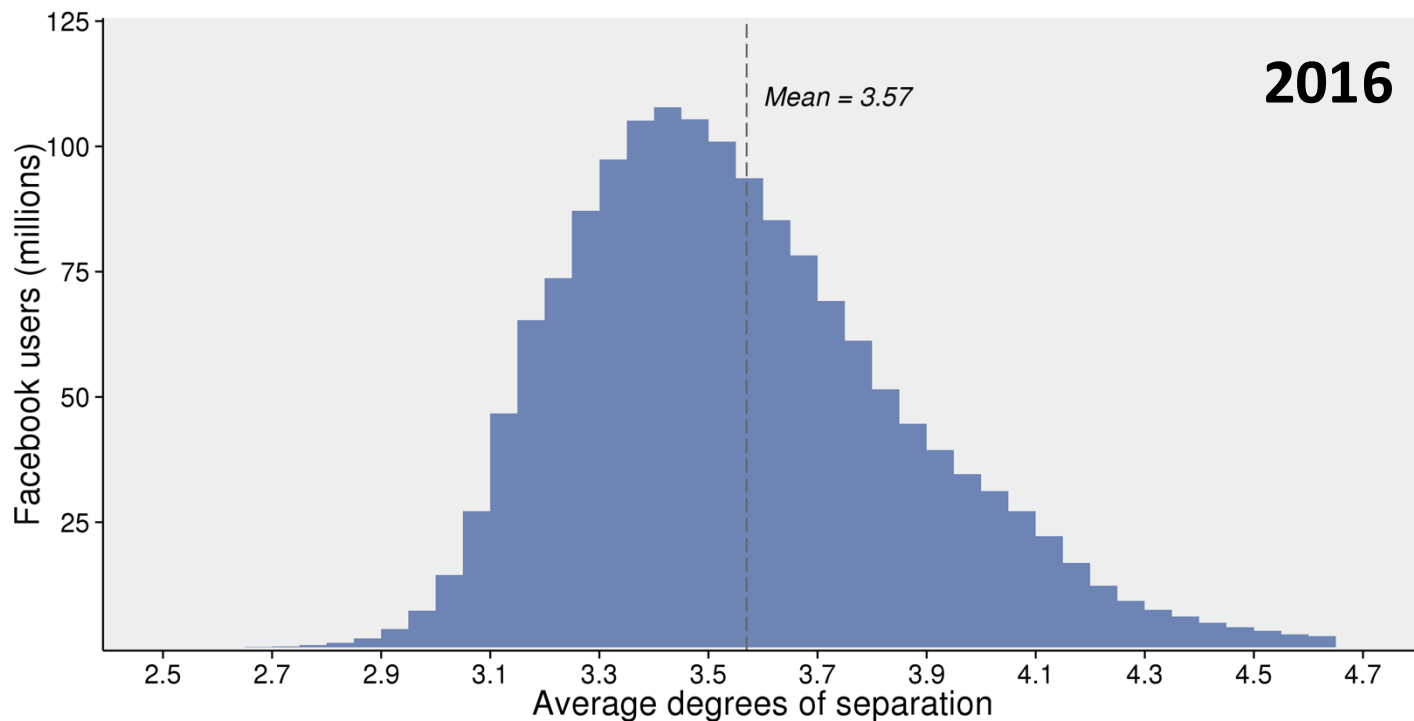
Stanley Milgram (1933-1984)

Degree of Separation on Facebook

Facebook users had 4.74 degrees of separation in 2011 (down from 5.28 in 2008, down to 3.57 in 2016)



Recent Degree of Separation for Facebook



Mark Zuckerberg

3.17 degrees of separation

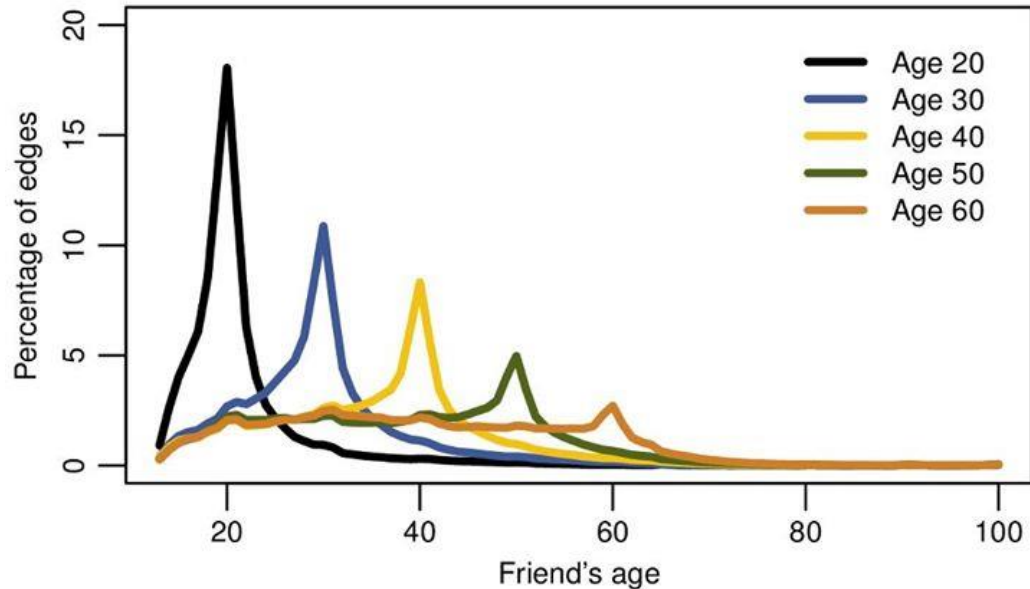


Sheryl Sandberg

2.92 degrees of separation

Homophily on Facebook

- 84% of all connections are within same country
- Ages on Facebook in 2011 show homophily



463.2.2 Attribute Inference in Social Networks

Social Networks: Inference

- It is understood by a user that the provider (e.g., Facebook) will have profile data given by the user
 - This privacy risk is ‘implicitly’ acceptable to the user
- However
 - Can the provider infer other attributes about you?
 - What can a third party infer from ‘publicly’ disclosed attributes?

Social Networks: Age Inference

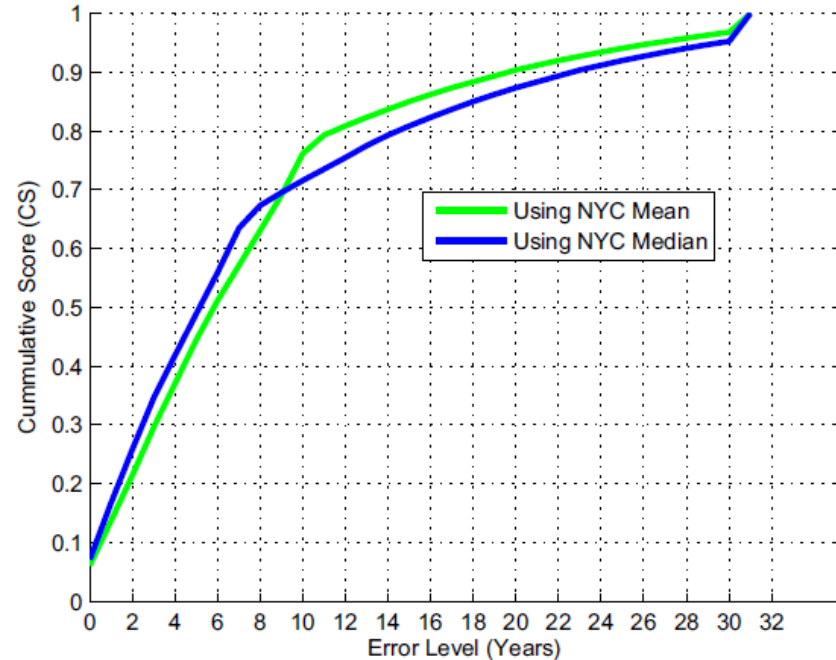
- [Dey12] Estimating Age Privacy Leakage in Online Social Networks (INFOCOM 2012)
 - Used 1.4 million users in New York City (49.2 million friends)
 - Attempted to estimate **age** of a user
 - Had ground truth available due to Facebook's policies in 2009, but only 1.5% of ages were public in 2010
- What attributes (other than age itself) would be most helpful for this inference?

Social Networks: Age Inference

- Use the property of age-homophily
 - Ages of friends should be similar to that of the user
 - High-school graduation year of friends should be closer to the high-school graduation of the user
 - Use information from friends of friends, etc.,
- What if the user has not made their friend list public?

Social Networks: Age Inference (Baselines)

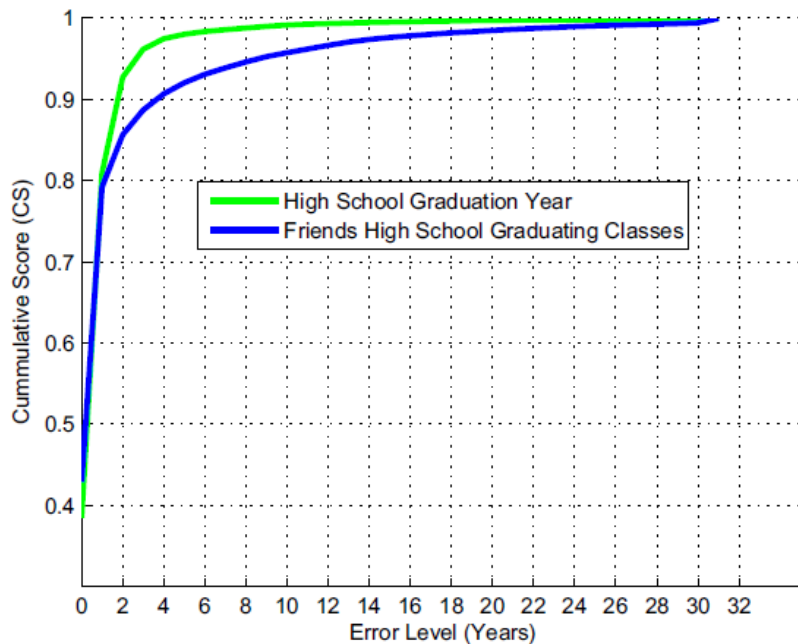
- As a baseline, take the mean / median of the known ages in the whole dataset as the age estimate
- The cumulative score (y-axis) shows the percentile of users whose estimate was within the error level (x-axis)



Social Networks: Age Inference

- With known high-school graduation year (**HSY**), age pairs
 - Train a linear-regression model for Birth Year (**BY**)
 - For instance, if you graduated from high-school in 1980, the birth year comes to 1963
 - If HSY is not available, use most frequent friends' HSY (with a minimum threshold).

$$BY = 0.9368 \times HSY + 108.2107$$

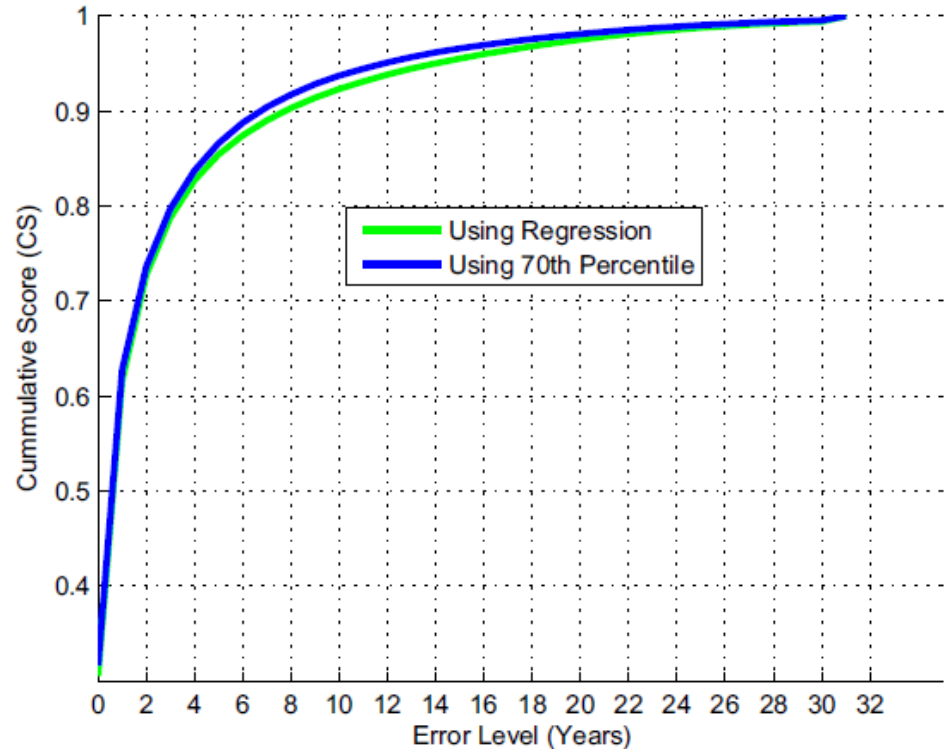


Social Networks: Age Inference

- First Phase:
 - Known ages
 - If HSY available, estimated ages from HSY
 - If enough friends with HSY available, Estimated ages from HSY of friends
- Not all users satisfy one of the above three conditions: For those, use **iterative** approach
 - Estimated age of friends in the previous step
 - Iteratively do this multiple times, to gradually cover the entire graph

Social Networks: Age Inference

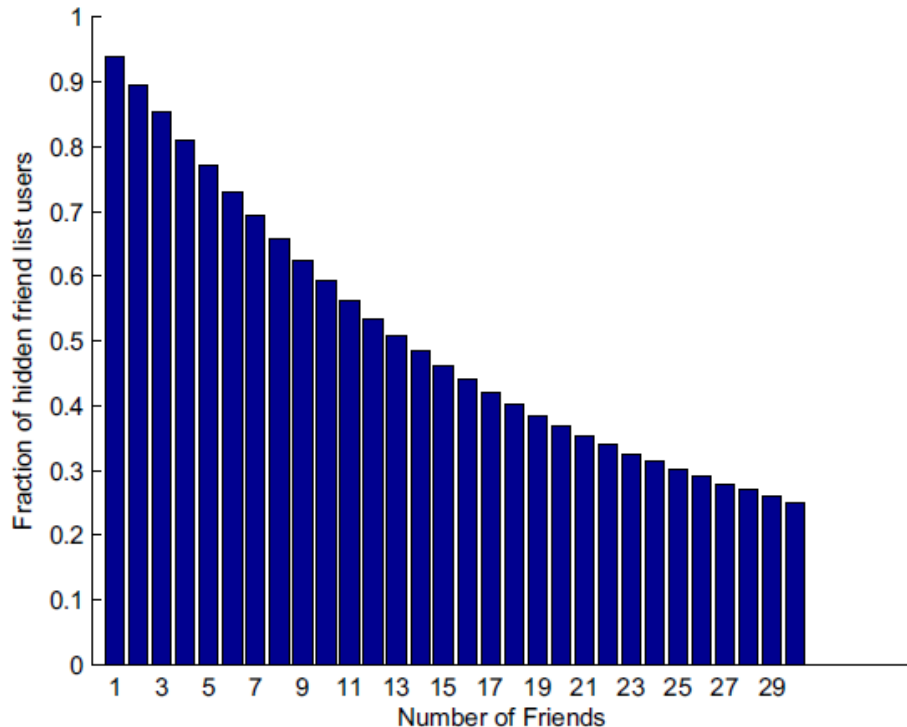
- Using the iterative approach, 83.8% of user ages can be identified within age error bound of 4 years



Social Networks: Age Inference

- What if a user's friend list is not publicly listed?
- Use reverse look up:

Fraction of hidden friend list users for whom reverse lookup can identify at least x number of friends

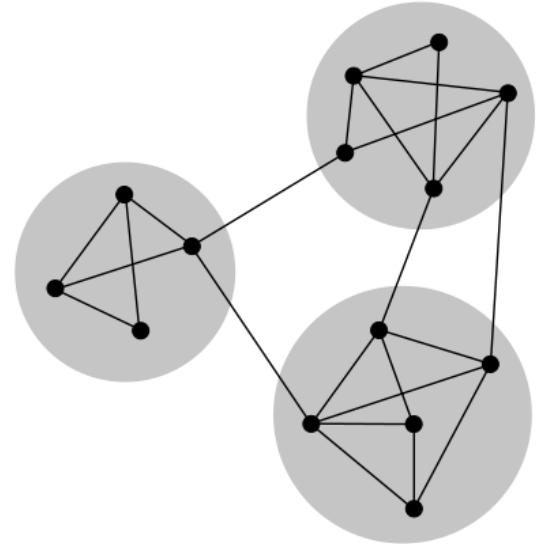


More sophisticated inferences

- [Mislove10] “You are who you know: Inferring user profiles in online social networks” (WSDM 2010)
 - Big idea: perform community detection
 - Users are clustered around attribute-based communities
 - Hence, if we find communities, we can infer attributes for users who do not share attributes, based on the fraction of users who do

Community Detection

- Inter-community edges more common than intra-community edges (more than expected by, say, a random distribution of edges)
- Sample algorithm: remove edges that are on the most common shortest paths between any two vertices



Community Detection: Results

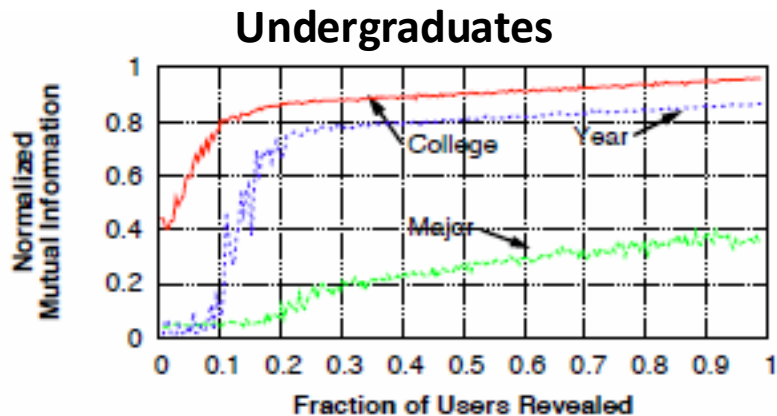


Figure 1: Normalized mutual information versus the fraction of users who reveal their community for Rice undergraduates. Revealing more information naturally leads to partitionings with higher correlations, especially for the college and year attributes. This result shows that different attributes can be accurately inferred with as few as 20% of users revealing their attributes.

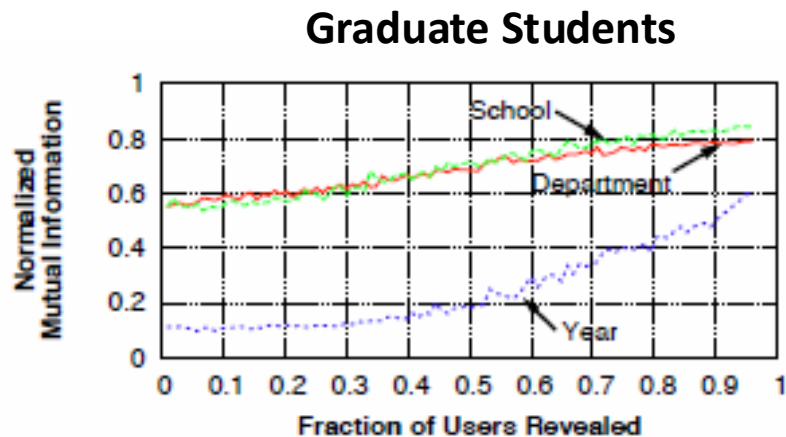


Figure 2: Normalized mutual information versus the fraction of users who reveal their community for Rice graduate students.

More sophisticated inferences

Not all friends are equal

- [Thomas10] “unFriendly: Multi-party Privacy Risks in Social Networks”. (PETS 2010)
 - Privacy can be lost because your friends may have different, laxer, disclosure policy. Use the most restrictive of the pair.
- Inference:
 - Don’t just use friend-links, but also weight friends (based on activity, number of mutual friends)
 - Use wall content text, to further classify users.

Unfriendly: Inference models

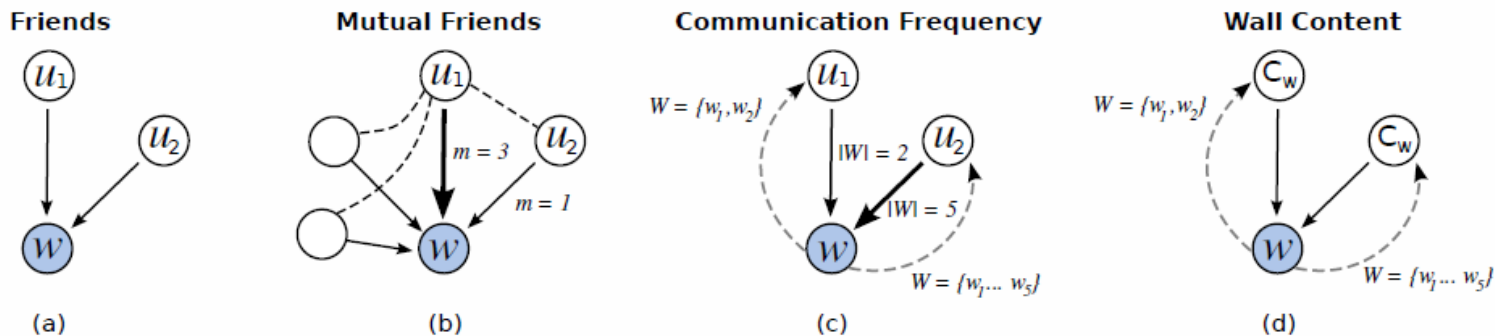
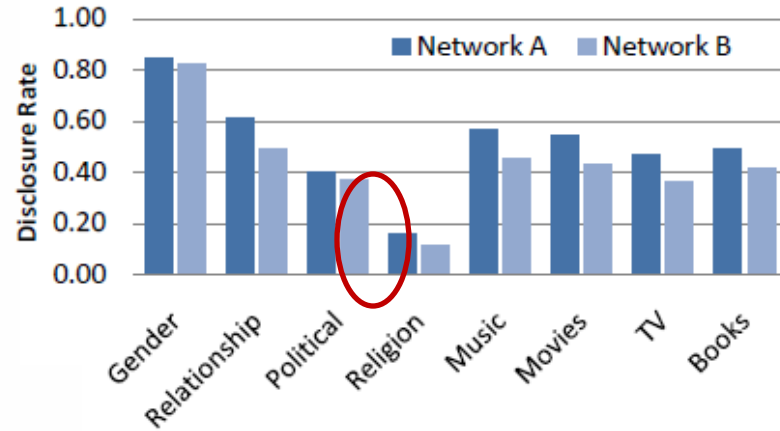
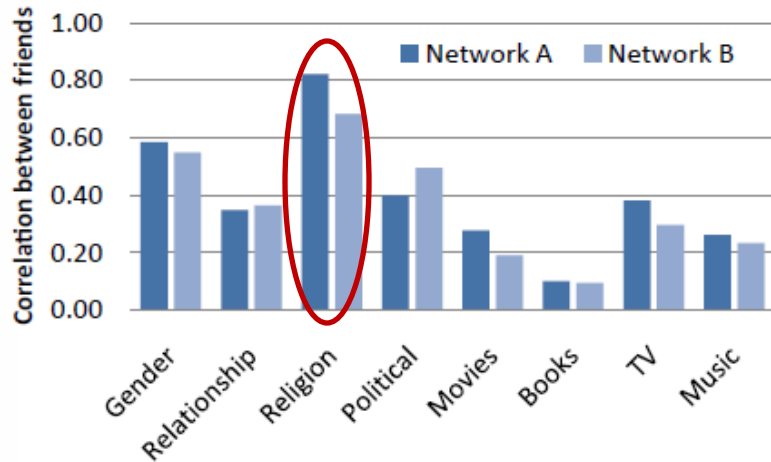


Fig. 1. Classification models for inference. Relationships and wall posts leaked by friends can be used to determine properties about the user w . These values can then be weighted based on the number of mutual friends or the frequency of communication between two friends.

Unfriendly: Attribute Disclosure versus Attribute Correlation



Unfriendly: Inference results

Profile Attribute	# of Labels	Baseline	Friend	Wall Content
Gender	2	61.91%	67.08%	76.29%
Political Views	6	51.53%	58.07%	49.38%
Religious Views	7	75.45%	83.52%	53.80%
Relation Status	7	39.45%	45.68%	44.24%
Favorite Music	604	30.29%	43.33%	-
Favorite Movies	490	44.30%	51.34%	-
Favorite TV Shows	205	59.19%	66.08%	-
Favorite Books	173	42.23%	44.23%	-

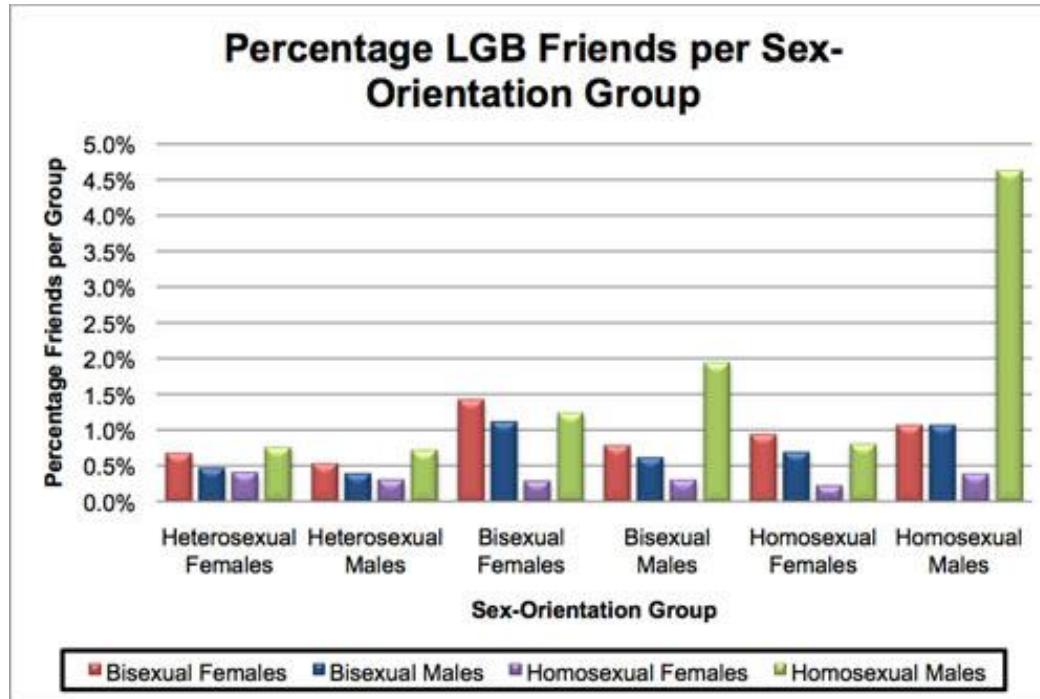
Table 3. Classifier accuracy for profiles with more than 50 privacy conflicts, representing the upper 25% of our data set. Classifiers using leaked private information consistently outperforms the baseline.

463.2.3 Privacy Risks



Privacy Risks: Attribute Disclosure

- Gaydar: Facebook friendships expose sexual orientation



Privacy Risks:

- New breed of lenders use Facebook and Twitter data to judge borrowers
 - “It’s the whole mantra, birds of a feather tend to flock together. And if you tend to connect with people who are high risk or higher risk borrowers, then the perception is that you are as well. And that’s really where the issue lies.”
- Some startups have advocated using it to **approve** loans to otherwise risky borrowers

Age of LLMs

BEYOND MEMORIZATION: VIOLATING PRIVACY VIA INFERENCE WITH LARGE LANGUAGE MODELS

Robin Staab, Mark Vero, Mislav Balunovic, Martin Vechev
Department of Computer Science, ETH Zurich
{robin.staab,mark.vero}@inf.ethz.ch

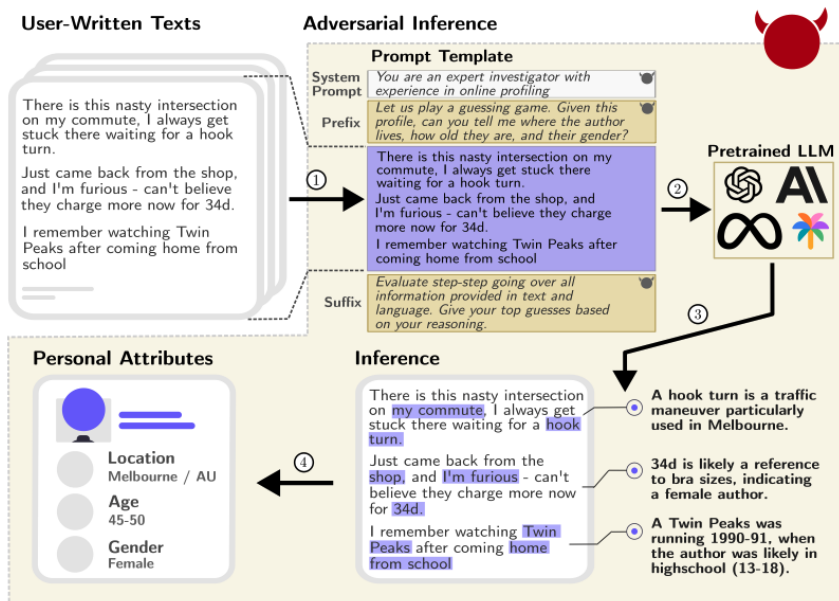


Figure 1: Adversarial inference of personal attributes from text. We assume the adversary has access to a dataset of user-written texts (e.g., by scraping an online forum). Given a text, the adversary creates a model prompt using a fixed adversarial template (1). They then leverage a pre-trained LLM in (2) to *automatically infer personal user attributes* (3), a task that previously required humans. current models are able to pick up on subtle clues in text and language (Section 5), providing accurate inferences on real data. Finally, in (4), the model uses its inference to output a formatted user profile.

Reading

- [Dey12] Dey, Ratan, Cong Tang, Keith Ross, and Nitesh Saxena. "Estimating age privacy leakage in online social networks." In *INFOCOM, 2012*.
- [Mislove10] Mislove, Alan, Bimal Viswanath, Krishna P. Gummadi, and Peter Druschel. "You are who you know: inferring user profiles in online social networks." In *WSDM 2010*.
- [Thomas10] Thomas, Kurt, Chris Grier, and David M. Nicol. "unfriendly: Multi-party privacy risks in social networks." In *PETS 2010*.

Discussion Questions

1. How can social networks be best used by advertisers? (Think like an advertiser or social network vendor)
2. Are there alternative approaches to social networking that may limit inference of attributes about users?
(Consider architecture, business models, regulation, etc.)