

Project presentation

Machine Learning and Big Data Processing(ELEC-Y591)

Cédric HANNOTIER Mathieu PETITJEAN Hasan Can YILDIRIM

June 13, 2018

Outline

- 1 State of the art
- 2 Approach
- 3 Results

State of the art

Outline

- 1 State of the art
- 2 Approach
- 3 Results

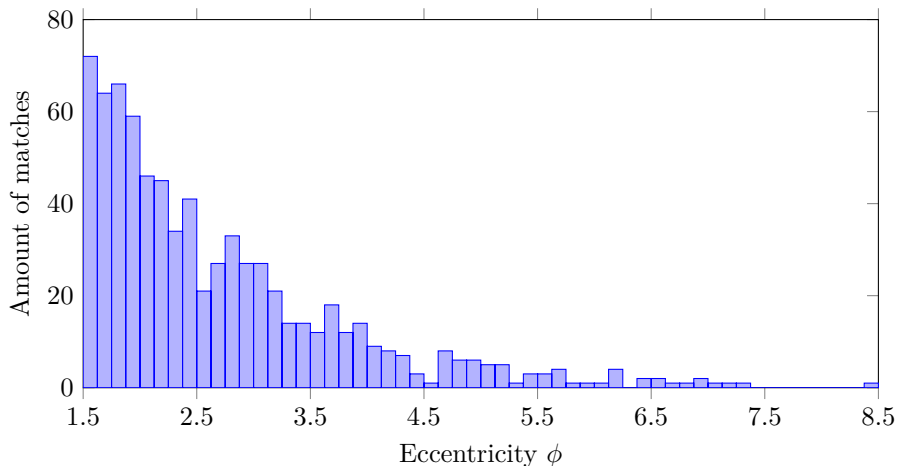
Approach

Outline

- 1 State of the art
- 2 Approach
- 3 Results**

Matches

- 3600 Netflix and 3000 MovieLens users
- a few more than 800 matches, 5 with $\phi > 7$

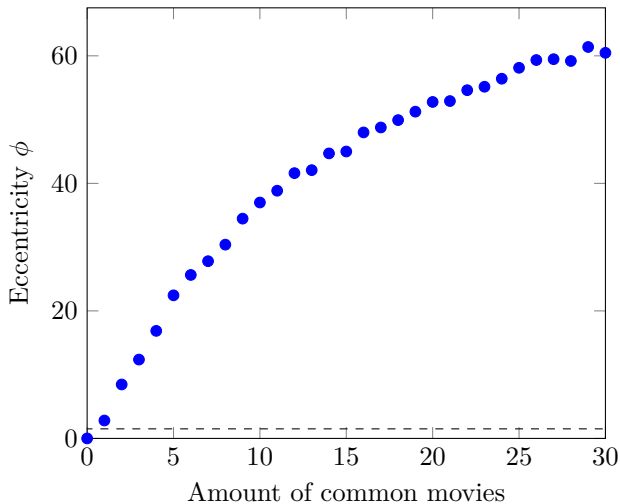


How to validate the results?

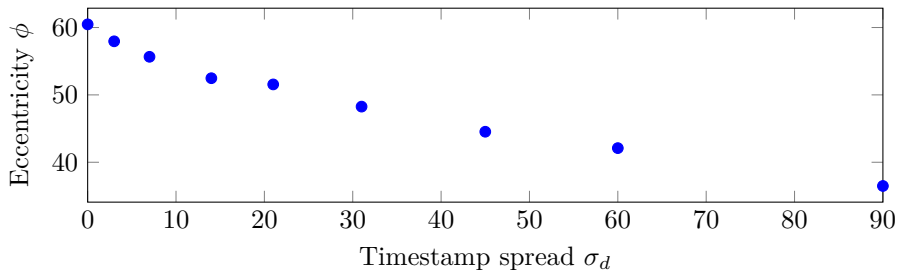
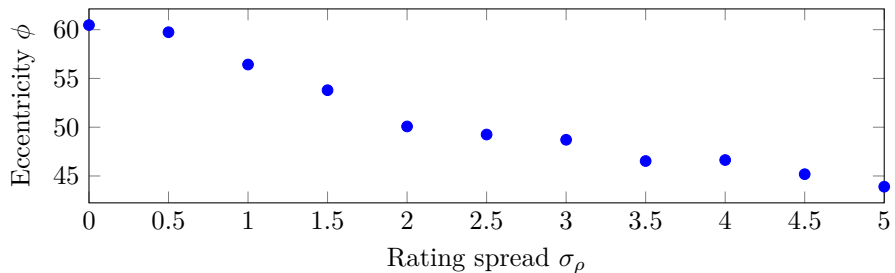
- No knowledge of the ground truth of matching users

Validation procedure

- Put a dummy user in both datasets
- Vary the number of common movies
- Perturb the rating and the timestamp with uniform distributed noise



Robustness



Conclusion

- $\phi_{\max} \approx 8.5$ while the validation process showed $\phi > 30$
- The original Netflix attack: $\phi_{\text{match}} = \{18, 25\}$

⇒ cannot conclude statistical quasi-certainty of de-anonymization

Possible improvements

- Use more than 0.04 % of the possible user combinations
- Tuning of the parameters ϕ, d_0, ρ_0
- Add more features (e.g. movie genres)
- Increase the errors impact on the scoring