

COMP 4332 / RMBI 4310

Big Data Mining (Spring 2023)

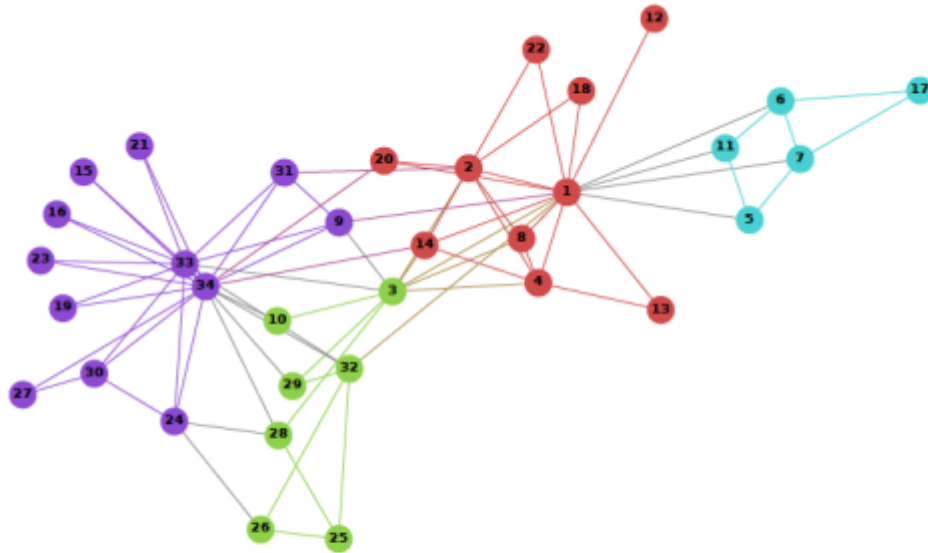
Project 2 Social Network Mining

TA: Weiqi Wang (wwangbw@connect.ust.hk)

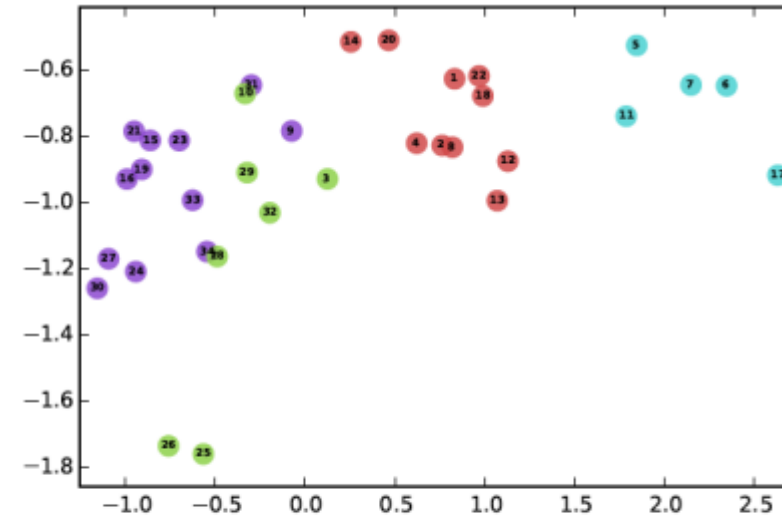
Social Network



Network Representations



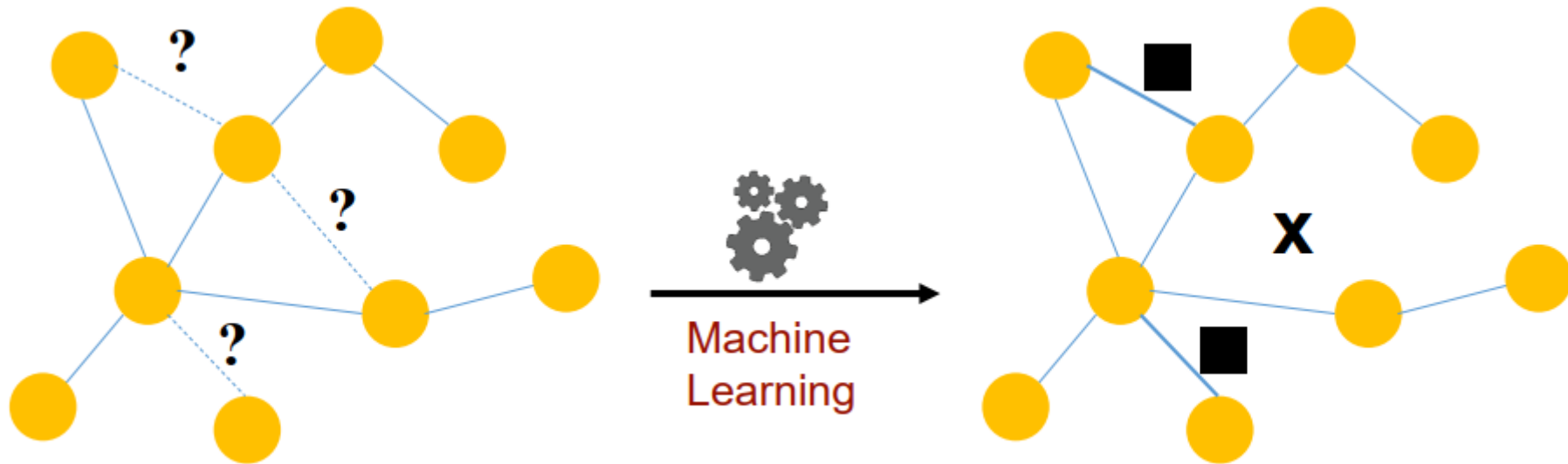
Input



Output

Link Prediction

- Predict the relation between nodes with their similarity and calculate the AUC-ROC score.



Pipeline

- Dataloader
- Random walk generator (first-order, second-order, ...)
- Embedding algorithm (DeepWalk, node2vec, ...)
- Scorer

Dataset

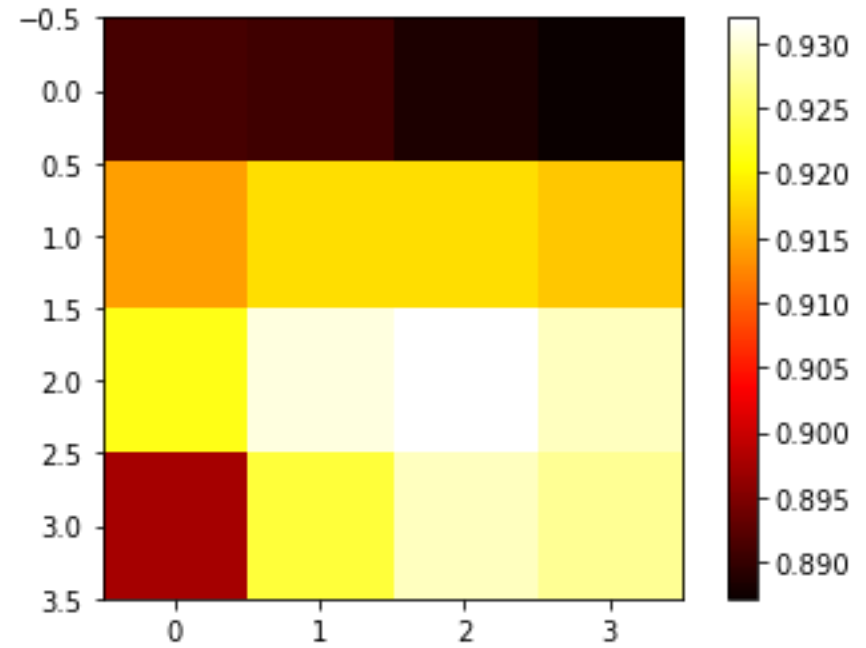
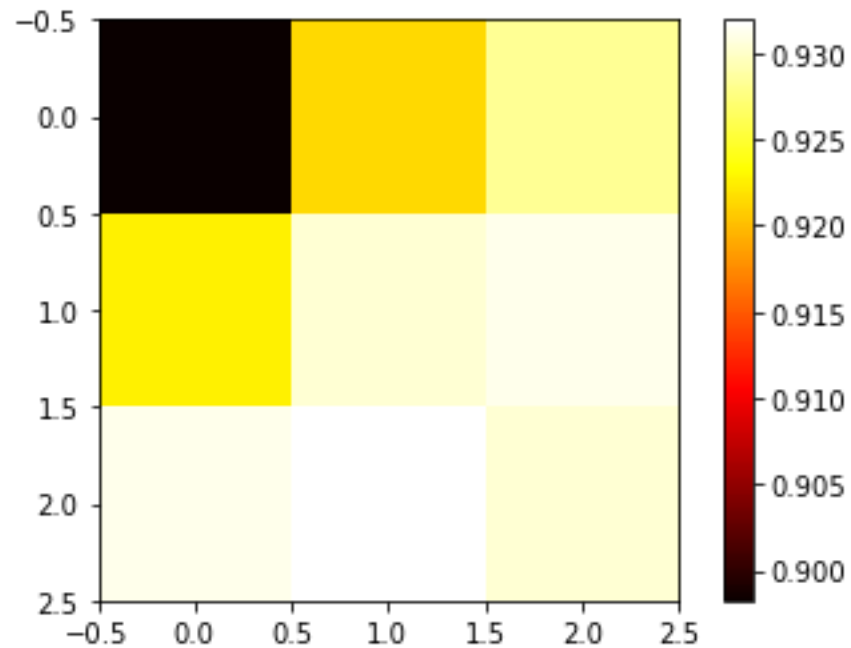
- training data: 8328 nodes, 100000 edges
- validation data: 5440 nodes, 19268 edges
- test data: 5452 nodes, 40000 edges (19267 positive edges)
- score: [0, 1]
- given features: `user_id`, `friends`

user_id	friends
--UOvCH5qEgdNQ8lzR8QYQ	[DBHCFW3mSmmOEpONHVu1rQ, QPJJohtGqkMkaN0Gt3TRI...
-05T0q5BxB9g0RCKiGYoyQ	[V7uS5US4oTf-S9u36HJQCQ]
-0HhZbPBIB1YZx3BhAfaEA	[uG35h72BAMutvXAWdRpqCQ, Sv48jgljDP-CRfXmU8uSg...
-1ZMRA0N01rqZL0TWk3fgA	[no2kFt4TEezZDVaM8haSDA]
-267Yx8RmdP6io2-qI4UcQ	[E3pXvQwKsPBQGQ7RkLrN3g, deL6e_z9xqZTIODKqnvRX...
...	...

src	dst	score
Nu5188fyBvHZHvgEgZT2bQ	IKSmm5MzHF8cMhMolKaIOW	0
rCWrxuRC8_pfagpchtHp6A	xQLy_wpqrR3etSxt61OIlg	0
maK3UBQczh33NuDjBYeHrA	K0sapHOIhIGNjx3GBesf5A	0
siXOnFrtV0a_YjOJr-X2Mg	7hAhYoMPjHnxKCz6MQ95Bg	0
_4iFWWuZ6_RrzyXZrMq3Mw	rTK_sTPgBjJXkdKi2G3X-w	0
...

Analysis

- Heatmap (<https://stackoverflow.com/questions/33282368/plotting-a-2d-heatmap-with-matplotlib>)



Evaluation

- AUC-ROC score on **test data**

Submission

- Predictions on **test data** (please make sure your pred.csv's format is same as test.csv: src/dst/score)
- Report (1~2 pages)
- Code (Frameworks and even programming languages are not restricted.)
- DDL: April 24, 2023
- Submission: Each **team leader** is required to submit the groupName.zip file that contains pred.csv, the report, and your team's code on canvas.
- We will check your report with your code and the AUC scores. You will be graded based on your testing set performance and your report.

Grading Rule

Grade	Model (80%)	Report (20%)	Baseline (on test set)
60%		submission	
80%	an easy baseline that most students can outperform	detailed explanation	40000 edges (20733 negative): 91.00%
90%	a competitive baseline that about half students can surpass	detailed explanation and analysis	40000 edges (20733 negative): 92.00%
100%	a very competitive baseline	excellent visualization and analysis	40000 edges (20733 negative): 93.00%

Thank You