# Data Analytics I| Assignment 3

**Author:** Prateek Pandey
**ID:** 2022201035

October 15, 2023



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

H Y D E R A B A D

# Contents

- After Joining Movies and Ratings Table

| | movieId | title | genres | userId | rating | timestamp |
|---|---|---|---|---|---|---|
| 0 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 964982703 |
| 1 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 847434962 |
| 2 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 7 | 4.5 | 1106635946 |
| 3 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 15 | 2.5 | 1510577970 |
| 4 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 17 | 4.5 | 1305696483 |
| ... | ... | ... | ... | ... | ... | ... |
| 100831 | 193581 | Black Butler: Book of the Atlantic (2017) | Action\|Animation\|Comedy\|Fantasy | 184 | 4.0 | 1537109082 |
| 100832 | 193583 | No Game No Life: Zero (2017) | Animation\|Comedy\|Fantasy | 184 | 3.5 | 1537109545 |
| 100833 | 193585 | Flint (2017) | Drama | 184 | 3.5 | 1537109805 |
| 100834 | 193587 | Bungo Stray Dogs: Dead Apple (2018) | Action\|Animation | 184 | 3.5 | 1537110021 |
| 100835 | 193609 | Andrew Dice Clay: Dice Rules (1991) | Comedy | 331 | 4.0 | 1537157606 |

100836 rows × 6 columns

- After Considering only those movies who have been rated > 2 and users who have rated more than 10 movies

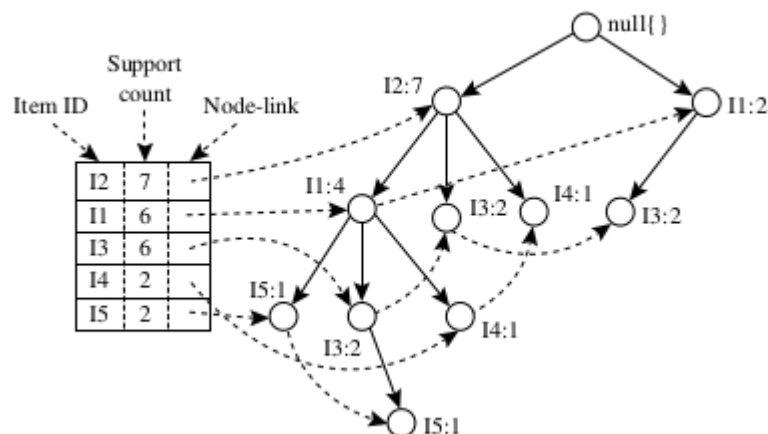| | userId | title |
|---|---|---|
| 0 | 1 | [Toy Story (1995), Grumpier Old Men (1995), He... |
| 1 | 2 | [Shawshank Redemption, The (1994), Tommy Boy (... |
| 2 | 3 | [Escape from L.A. (1996), Highlander (1986), S... |
| 3 | 4 | [Get Shorty (1995), To Die For (1995), Mighty ... |
| 4 | 5 | [Toy Story (1995), Get Shorty (1995), Babe (19... |
| ... | ... | ... |
| 605 | 606 | [Toy Story (1995), Sabrina (1995), American Pr... |
| 606 | 607 | [Toy Story (1995), American President, The (19... |
| 607 | 608 | [Toy Story (1995), GoldenEye (1995), Casino (1... |
| 608 | 609 | [Toy Story (1995), GoldenEye (1995), Bravehear... |
| 609 | 610 | [Toy Story (1995), Heat (1995), Casino (1995),... |

607 rows × 2 columns

- Dividing the data set into 80% training set and 20% test set for each user

| | userId | training_set | test_set |
|---|---|---|---|
| 0 | 1 | [Toy Story (1995), Grumpier Old Men (1995), He... | [Kiss the Girls (1997), Wild Things (1998), To... |
| 1 | 2 | [Shawshank Redemption, The (1994), Tommy Boy (... | [Inside Job (2010), Collateral (2004), Gladiat... |
| 2 | 3 | [Escape from L.A. (1996), Highlander (1986), S... | [Saturn 3 (1980), Clonus Horror, The (1979), C... |
| 3 | 4 | [To Die For (1995), Mighty Aphrodite (1995), P... | [Name of the Rose, The (Name der Rose, Der) (1... |
| 4 | 5 | [Toy Story (1995), Dead Man Walking (1995), Us... | [Lion King, The (1994), Clueless (1995), Babe ... |
| ... | ... | ... | ... |
| 605 | 606 | [Sabrina (1995), American President, The (1995... | [Romeo and Juliet (1968), What Have I Done to ... |
| 606 | 607 | [Toy Story (1995), American President, The (19... | [Saving Private Ryan (1998), Lady and the Tram... |
| 607 | 608 | [GoldenEye (1995), Casino (1995), Get Shorty (... | [Groundhog Day (1993), Addams Family, The (199... |
| 608 | 609 | [Toy Story (1995), GoldenEye (1995), Bravehear... | [Net, The (1995), Return of Martin Guerre, The... |
| 609 | 610 | [Toy Story (1995), Heat (1995), Casino (1995),... | [Goodnight Mommy (Ich seh ich seh) (2014), Pre... |

607 rows × 3 columns

- Using FP-Growth Algorithm we first created an FP-TREE

- After Mining Frequent Patterns on the FP-Tree resulting Frequent Patterns
  are

```
('Back to the Future (1985)', 'Fight Club (1999)'): 55,
('Back to the Future (1985)', 'Shrek (2001)'): 59,
('Back to the Future (1985)', 'Terminator 2: Judgment Day (1991)'): 62,
('Back to the Future (1985)', 'Silence of the Lambs, The (1991)'): 62,
('Back to the Future (1985)', 'Toy Story (1995)'): 64,
('Back to the Future (1985)', 'Jurassic Park (1993)'): 64,
('Back to the Future (1985)', 'Matrix, The (1999)'): 66,
('Back to the Future (1985)', 'Shawshank Redemption, The (1994)'): 66,
('Back to the Future (1985)',
 'Raiders of the Lost Ark (Indiana Jones and the Raiders of the Lost Ark) (1981)'): 68,
('Back to the Future (1985)',
 'Raiders of the Lost Ark (Indiana Jones and the Raiders of the Lost Ark) (1981)',
 'Star Wars: Episode V - The Empire Strikes Back (1980)'): 52,
('Back to the Future (1985)',
 'Star Wars: Episode IV - A New Hope (1977)',
 'Star Wars: Episode V - The Empire Strikes Back (1980)'): 50,
('Back to the Future (1985)',
 'Star Wars: Episode IV - A New Hope (1977)',
 'Star Wars: Episode VI - Return of the Jedi (1983)'): 55,
('Back to the Future (1985)', 'Pulp Fiction (1994)'): 70,
('Back to the Future (1985)',
 'Forrest Gump (1994)',
 'Pulp Fiction (1994)'): 51,
('Back to the Future (1985)',
 'Star Wars: Episode V - The Empire Strikes Back (1980)'): 72,
('Back to the Future (1985)',
 'Star Wars: Episode V - The Empire Strikes Back (1980)',
 'Star Wars: Episode VI - Return of the Jedi (1983)'): 54,
```

- Mining all Association Rules

```
[((('Star Wars: Episode II - Attack of the Clones (2002)',),
  ('Star Wars: Episode VI - Return of the Jedi (1983)',),
  0.8333333333333334,
  55),
 (('Sin City (2005)',), ('Fight Club (1999)',), 0.828125, 53),
 (('Crimson Tide (1995)',), ('Fugitive, The (1993)',), 0.75, 54),
 (('Sleepless in Seattle (1993)',),
  ('Forrest Gump (1994)',),
  0.7435897435897436,
  58),
 (('Heat (1995)',),
  ('Silence of the Lambs, The (1991)',),
  0.7162162162162162,
  53),
 (('Big Lebowski, The (1998)',),
  ('Pulp Fiction (1994)',),
  0.7142857142857143,
  50),
 (('O Brother, Where Art Thou? (2000)',),
  ('Matrix, The (1999)',),
  0.6986301369863014,
  51),
 (('Harry Potter and the Chamber of Secrets (2002)',),
  ('Harry Potter and the Prisoner of Azkaban (2004)',),
  0.6986301369863014,
  51),
 (('Ghost (1990)',), ('Forrest Gump (1994)',), 0.6973684210526315, 53),
 (('Requiem for a Dream (2000)',),
  ('Fight Club (1999)',),
  0.6944444444444444,
  50)]
```

- Recommendation-

  ○ After Mining Association Rules over the frequent patterns

    ▪ Sorting these rules based on confidence

- Sorting these rules based on Support

- Taking Intersection of both these top 100 sorted rules by support and confidence to get common list

- Precision and Recall Graph of Random 25 users

- Given rule X → Y, with X representing a movie from the training set, set Y comprises recommended movies. By aggregating the movies on the right side of N rules, denoted as R, the resulting set serves as recommendations. The overlap between R and the test set is termed the hit set. The recall is determined by the ratio of the hit set to the test set, while the precision is calculated as the ratio of the hit set to the recommendation set.
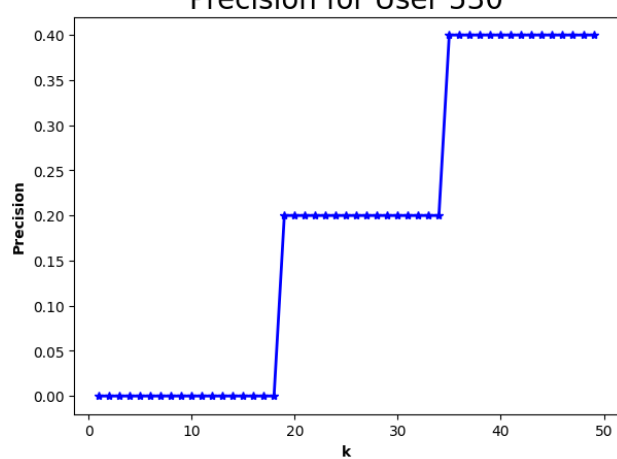
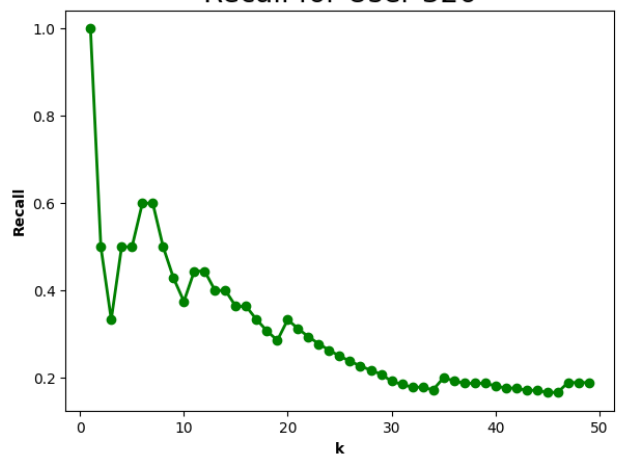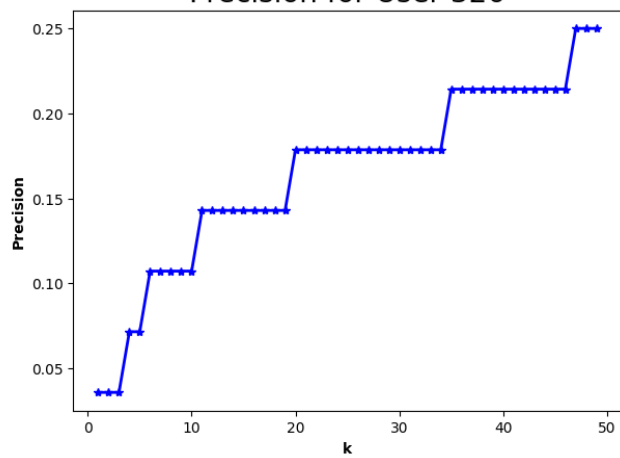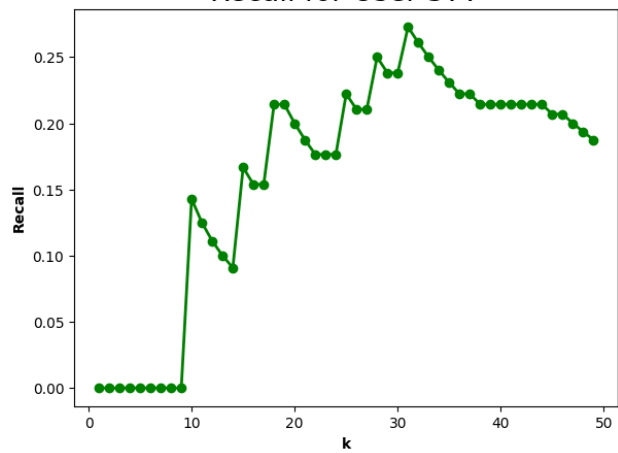Recall for User 325

Precision for User 325
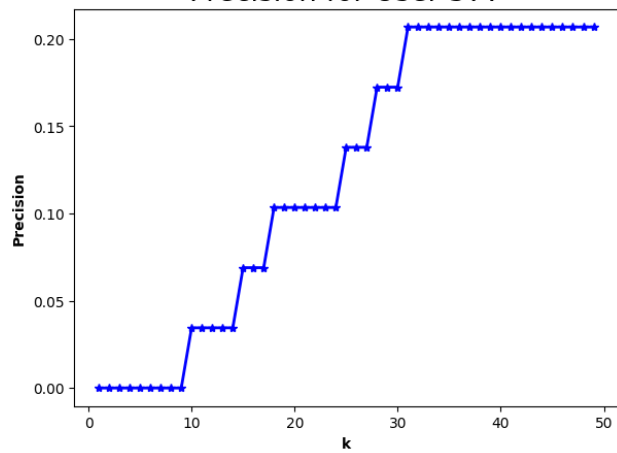
Recall for User 430

Precision for User 430
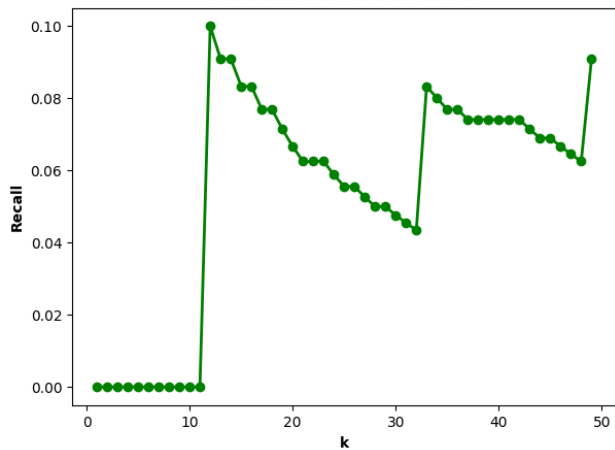
Recall for User 186
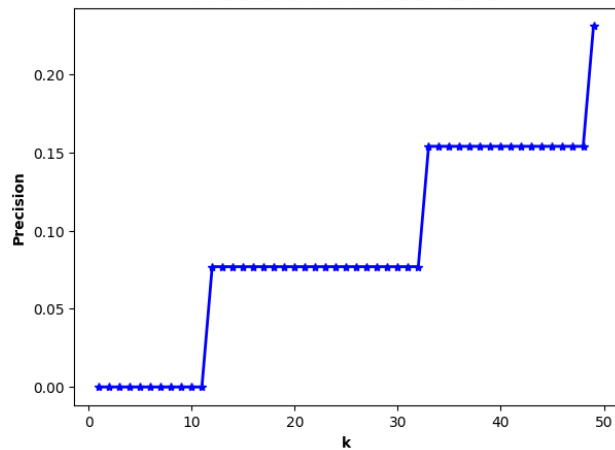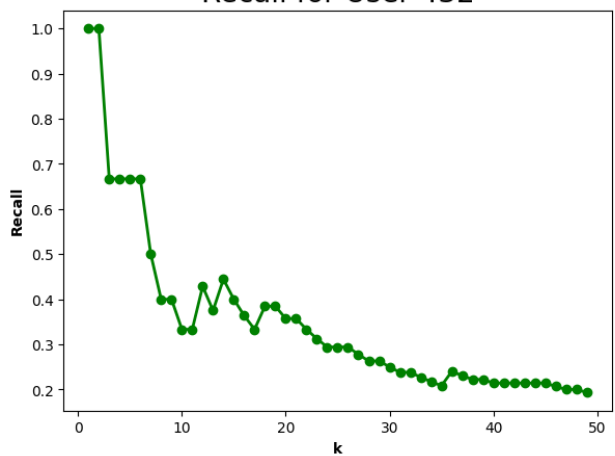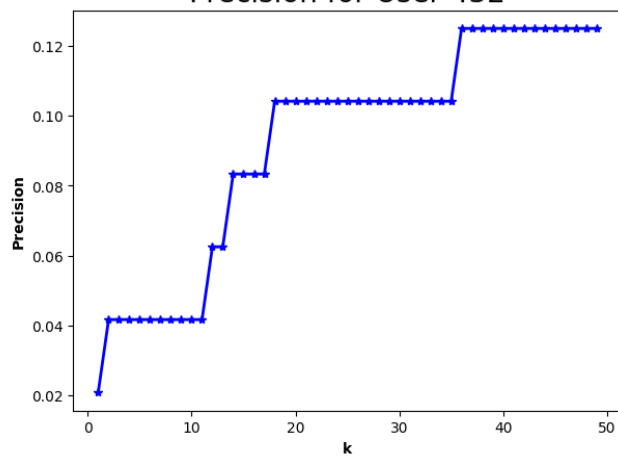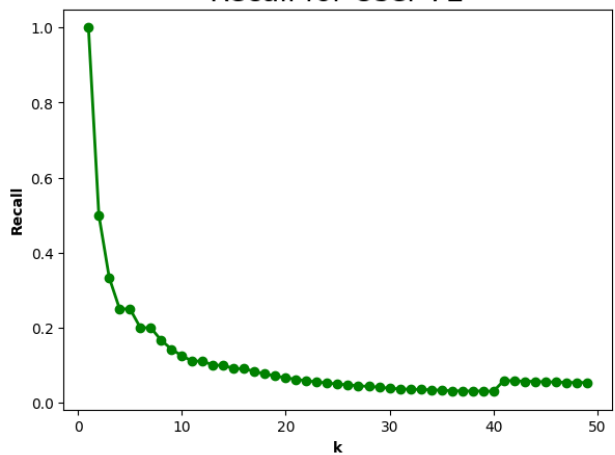
Precision for User 186

**Recall for User 473** · **Precision for User 473**

**Recall for User 14** · **Precision for User 14**

**Recall for User 8** · **Precision for User 8**