

Introduction:

This project is developed for the needs of investors in Indian Mutual funds market. It considers using association rule mining algorithm and collaborative filtering to provide personalized investment recommendations to users based on their risk profile, investment objectives, and financial preferences. This system assists users in making informed investment decisions in mutual funds. By analyzing transactional data of mutual funds and client attributes, the system aims to match users with the most suitable mutual fund options, thereby optimizing their investment portfolio and maximizing returns.

Scope:

- Collection and preprocessing of mutual fund data from the Value Research website.
- Generation of synthetic client data.
- Correlation analysis between mutual fund risk factors and client risk profiles.
- Creation of transactional dataset.
- Implementation of association rule mining and item-item collaborative filtering techniques for recommendation generation.

Dataset description:

The MF data is taken from Value research website, about Indian mutual funds' performance and attributes. The dataset includes funds from various categories, each represented by separate CSV files containing approximately 50 funds of a specific category.

Attributes of the mutual fund data include returns (1 wk, 3 m, 6 m, 1 yr, 3 yr, 5 yr, and 10 yr), market capitalization, turnover, net assets, riskometer, standard deviation, expense ratio, and other relevant metrics.

Client data is randomly generated and labeled based on current requirements. It contains attributes employment status, income range, age, investment horizon, investment objective, and risk appetite rating derived from psychological and market-based questions.

The recommendation process involves matching the risk profiles of MF with those of clients and allocating suitable funds to client baskets. The resulting transactional data is then binary encoded for further processing.

Data cleaning: Filling missing value in MF data.

Data encoding: Categorical to numeric transformation in case of MF categories, riskometer etc.

Data transformation: Binary encoded transactional data for applying asso. rule mining.

Novelty:

- **Multiple Data Sources:** We used data from MF and client both, with multiple attributes.
- **Recommendation Techniques and Analysis:** Apriori, FP growth and collaborative filtering, to generate personalized investment recommendation.
- **Risk Profile Matching:** A key novelty of our approach is the incorporation of risk profile matching between MF and clients.
- **Dynamic Recommendation Generation:** If new user or MF joins the data, using regression and classifier, values of derived attributes can be predicted.

Algorithm:

Algorithm 1 Mutual Fund Recommendation Algorithm

Require: Mutual fund data D_{funds} , Client data D_{clients}

Ensure: Recommended funds for each client

- 1: Calculate risk ρ_f for all mutual funds in D_{funds}
- 2: Calculate risk ρ_c for all clients in D_{clients}
- 3: **for** each client C_i in D_{clients} **do**
- 4: **for** each mutual fund MF_j in D_{funds} **do**
- 5: Calculate compatibility between ρ_f and ρ_c
- 6: **end for**
- 7: Select top 5 MF for C_i
- 8: **end for**
- 9: Transaction data $\xrightarrow{\text{convert}}$ Binary encoded data
- 10: Apply Apriori
- 11: Apply FP-Growth
- 12: Time analysis of Apriori and FP-Growth
- 13: Correlation analysis and Chi-square analysis on the association rules

Algorithm 2 Item-Item Collaborative Filtering

Require: Mutual fund data D_{funds} , Client data D_{clients}

Ensure: Recommended funds for each client

- 1: Calculate risk ρ_f for all mutual funds in D_{funds}
- 2: Calculate risk ρ_c for all clients in D_{clients}
- 3: Generate user-item matrix M based on risk matching as in step 3,4,5 in Algo.1
- 4: **for** each client C_i in D_{clients} **do**
- 5: Extract client investments from M
- 6: Predicted ratings using item-item collaborative filtering
- 7: Select top 5 MF for client C_i
- 8: **end for**
- 9: Analysis and Inference

Algorithm 3 PrefixSpan Algorithm

Require: Sequences S , Minimum support threshold $min_support$

Ensure: Frequent sequential patterns $freq_patterns$

- 1: $freq_patterns \leftarrow \{\}$
PrefixSpanprojected_db, prefix, min_support
 - 2: $item_count \leftarrow \{\}$
 - 3: **for** each sequence seq in $projected_db$ **do**
 - 4: **for** each item $item$ in seq **do**
 - 5: **if** $item$ exists in $item_count$ **then**
 - 6: Increment count of $item$ in $item_count$
 - 7: **else**
 - 8: Add $item$ to $item_count$ with count 1
 - 9: **end if**
 - 10: **end for**
 - 11: **end for**
 - 12: **for** each item $item$, count $count$ in $item_count$ **do**
 - 13: **if** $count \geq min_support$ **then**
 - 14: $new_prefix \leftarrow prefix \cup \{item\}$
 - 15: Add new_prefix to $freq_patterns$
 - 16: $projected_db_new \leftarrow \{seq[seq.index(item) + 1 :] \text{ for } seq \text{ in } projected_db \text{ if } item \text{ in } seq\}$
 - 17: **if** $projected_db_new$ is not empty **then**
 - 18: PREFIXSPAN($projected_db_new$, new_prefix , $min_support$)
 - 19: **end if**
 - 20: **end if**
 - 21: **end for**
 - 22: PREFIXSPAN(S , $\{\}$, $|S| \times min_support$)
-

Results and Discussion:

Time Analysis:

FP-Growth generally outperformed Apriori in terms of execution time, especially for larger datasets and lower support thresholds. This can be attributed to the efficient data structure (tree) used by the FP-Growth algorithm, which reduces the computational overhead associated with candidate generation and pruning.

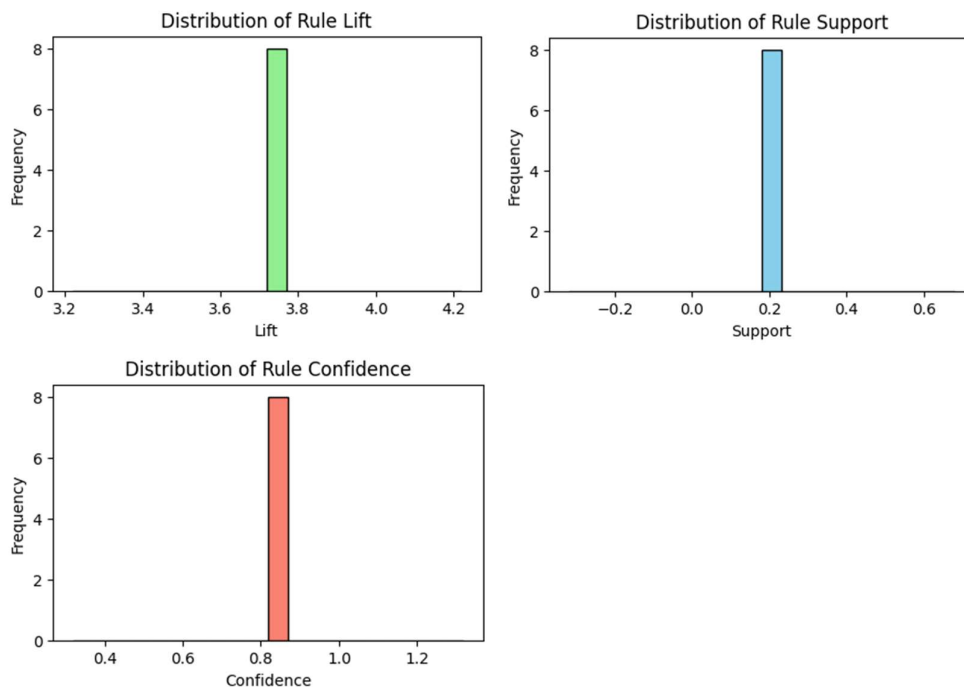
Correlation analysis:

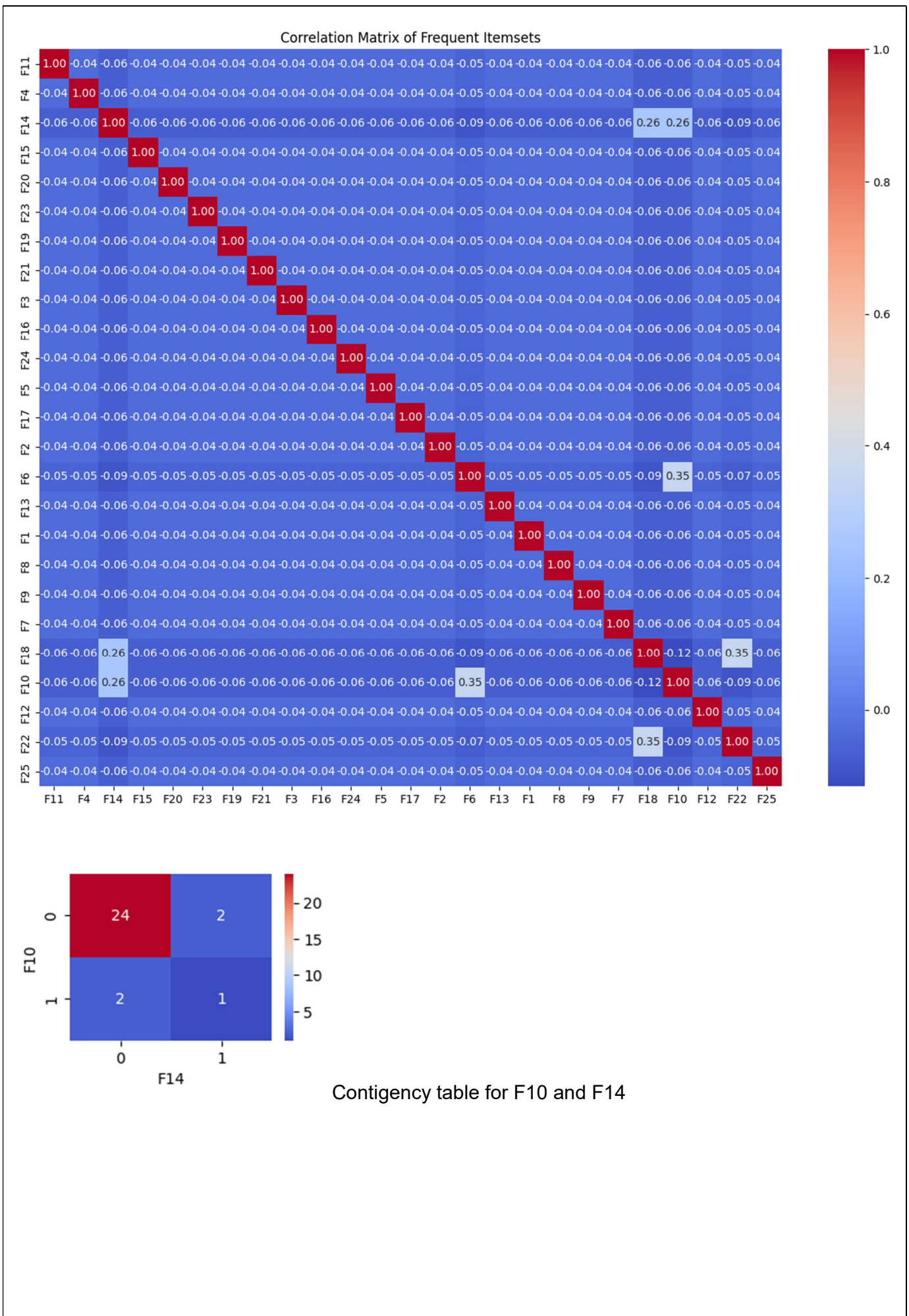
Correlation coefficients were computed for pairs of mutual funds based on their co-occurrence in transactions. The results were visualized using a heatmap, highlighting strong positive and negative correlations between funds.

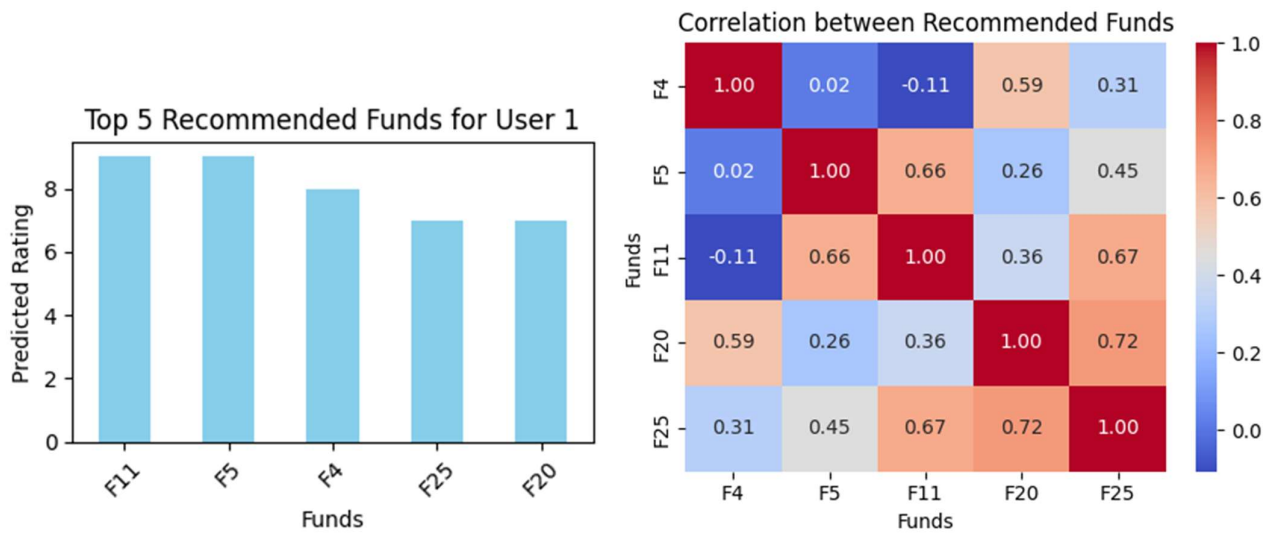
Chi-Square analysis:

It was performed to assess the statistical significance of association rules. Contingency tables formed by the presence and absence of mutual funds in transactions. The resulting p-values were used to determine the strength of association between funds, with lower p-values indicating higher significance.

Plots:







RMSE for predicted ratings:

```
rmse = calculate_rmse(df_nan, ratings)
print("RMSE:", rmse)
```

RMSE: 3.866908285771863

Overall, the results demonstrate the effectiveness of the Mutual Fund Recommendation System in generating meaningful investment recommendations based on association rule mining techniques. Further refinement and optimization of the system can be DONE to enhance its accuracy.