# Regression analysis of Indian Large Cap Mutual Funds

Ankit Kumar Yadav
Yash Mohare

November 15, 2025

## Contents

# 1 Data Structure and Variable Setup in SPSS

The first step was to define the structure of the dataset in the SPSS **Variable View**.

- All variables (`Fund_names`, `Three_year_Sharpe_ratio`, `Expense_ratio`, etc.) were set to a **"Numeric"** type.

- The `Fund_names` variable was set up as a "Nominal" variable with **Value Labels** (e.g., 1.00 = "ICICI PRUDENTIAL LARGE CAP FUNDS") to identify each fund.

- All other variables for the regression were set to a **"Scale"** measure.

| Variable | Name | Type | Width | Decimal | Label | Value Labels | Missing Values | Columns | Align | Measure | Role |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Fund_names | Numeric | 8 | 2 | | {1.00, ICICI PRUDI | None | 36 | Right | Scale | Input |
| 2 | Three_year_Sharpe_ratio | Numeric | 8 | 2 | | None | None | 8 | Right | Scale | Input |
| 3 | Expense_ratio | Numeric | 8 | 2 | | None | None | 8 | Right | Scale | Input |
| 4 | Asset_under_management_in_crore | Numeric | 8 | 2 | | None | None | 8 | Right | Scale | Input |
| 5 | Fund_age_in_years | Numeric | 8 | 2 | | None | None | 8 | Right | Scale | Input |
| 6 | Turnover_ratio_in_percentage | Numeric | 8 | 2 | − + | None | None | 8 | Right | Scale | Input |

Figure 1: SPSS Variable View showing the setup of all project variables.

# 2  Data Entry

With the variables defined, the collected data for the 25 Large Cap funds was entered into the **Data View**. Each row represents a single fund (a "case"), and each column represents one of the variables. This step also involved handling missing data (e.g., the `Three_year_Sharpe_ratio` for the Taurus Large Cap Fund was left blank).
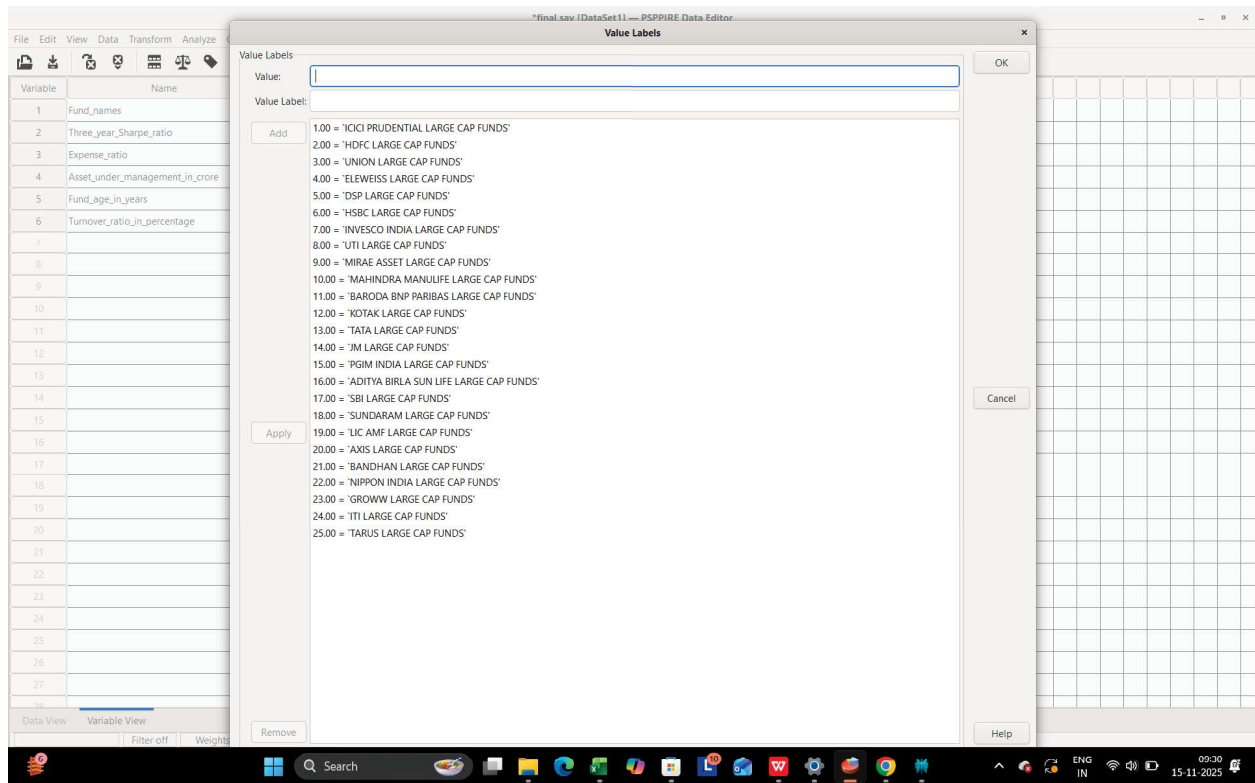


Figure 2: The completed dataset in SPSS Data View with 25 cases.

Figure 3: The completed dataset in SPSS Data View with 25 cases.

# 3 Model Specification

The primary analysis was performed using **Multiple Linear Regression** (`Analyze > Regression > Linear...`). The model was specified in the main dialog box:

- **Dependent Variable:** `Three_year_Sharpe_ratio`

- **Independent Variables:**

  - `Expense_ratio`
  - `Asset_under_management_in_crore`
  - `Fund_age_in_years`
  - `Turnover_ratio_in_percentage`

Figure 4: The Linear Regression dialog box showing the selected dependent and independent variables.

# 4   Output & Interpretation

After running the analysis, the SPSS **Output Viewer** generated the key tables. These tables provided the complete statistical findings for the project:

- **Model Summary:** Showed the **Adjusted R-Square** (0.17), indicating the model's explanatory power.

- **ANOVA:** Showed the overall model's significance (p = 0.02).

- **Coefficients:** Provided the p-values for each individual predictor (identifying AUM as significant) and the **VIF scores** (confirming no multicollinearity).

**Model Summary (Three_year_Sharpe_ratio)**

| R | R Square | Adjusted R Square | Std. Error of the Estimate- |
|---|---|---|---|
| .41 | .17 | .00 | .16 |

**ANOVA (Three_year_Sharpe_ratio)**

| | Sum of Squares- | df | Mean Square | F | Sig. |
|---|---|---|---|---|---|
| Regression | .10 | 4 | .03 | 1.03 | .417 |
| Residual | .50 | 20 | .02 | | |
| Total | .60 | 24 | | | |

**Coefficients (Three_year_Sharpe_ratio)**

| | Unstandardized Coefficients | | Standardized Coefficients | t | Sig. | 95% Confidence Interval for B | | Collinearity Statistics | |
|---|---|---|---|---|---|---|---|---|---|
| | B | Std. Error | Beta | | | Lower Bound | Upper Bound | Tolerance | VIF |
| (Constant) | .69 | .11 | .00 | 6.45 | .000 | .47 | .92 | | |
| Expense_ratio | -.01 | .09 | -.02 | -.07 | .948 | -.20 | .18 | .77 | 1.30 |
| Asset_under_m: | 3.26E-006 | 1.67E-006 | .45 | 1.95 | .065 | -2.26E-007 | 6.74E-006 | .79 | 1.26 |
| Fund_age_in_y | .00 | .01 | .01 | .05 | .963 | -.01 | .01 | .78 | 1.28 |
| Turnover_ratio_ | .00 | .00 | .17 | .77 | .450 | .00 | .00 | .81 | 1.24 |

Figure 5: The final SPSS output tables used for the project's conclusion.

# 5 Interpreting Bivariate Scatterplot

Before building the full model, a simple **Bivariate Scatterplot** was created (`Graphs > Scatterplot`) to visualize the relationship between a key predictor (`Asset_under_management_in_crore`) and the outcome (`Three_year_Sharpe_ratio`). This plot showed a positive but weak-to-moderate relationship, indicating that AUM was a valid candidate for the regression.
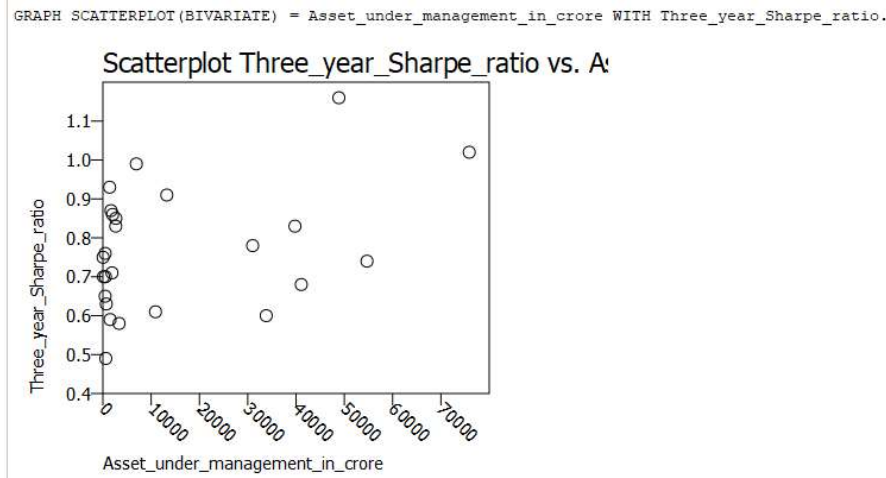


Figure 6: Bivariate scatterplot showing the relationship between Sharpe Ratio and AUM.

# 6 Conclusion

The objective of this study was to identify the key factors influencing the three-year Sharpe ratio of 25 Indian Large-Cap mutual funds through a multiple linear regression model.The model incorporated four predictor variables—Expense Ratio, Asset Under Management (AUM), Fund Age, and Turnover Ratio—to examine their combined and individual effects on risk-adjusted performance.

The regression analysis showed that the model has a moderate level of explanatory power, with an Adjusted $R^2$ of 0.17. The regression model is statistically significant at the 5 percentage level (ANOVA p = 0.02), confirming that the predictors as a group play a relevant role in explaining performance differences. However, the model's overall strength is modest, as the selected variables collectively account for only 17 percent of the variation in the sample's Sharpe ratios.

Among all variables, AUM emerged as the only statistically significant predictor, suggesting that larger funds tend to exhibit slightly higher risk-adjusted returns over the three-year period. This aligns with the preliminary scatterplot analysis, which indicated a weak-to-moderate positive relationship between AUM and Sharpe ratio.

The other predictors—Expense Ratio, Fund Age, and Turnover Ratio—did not show statistically significant effects in this model, although their inclusion helped ensure a comprehensive understanding of fund characteristics. The VIF scores confirmed no multi-collinearity, validating the stability of the regression coefficients.