# Data Collection and Preprocessing Phase

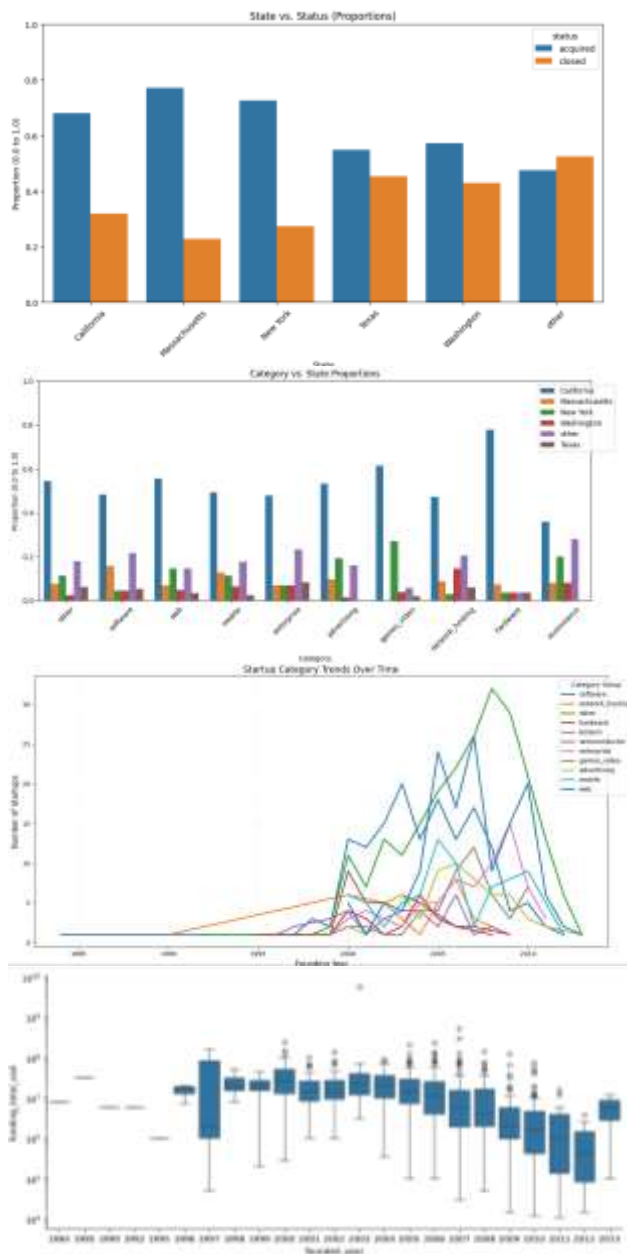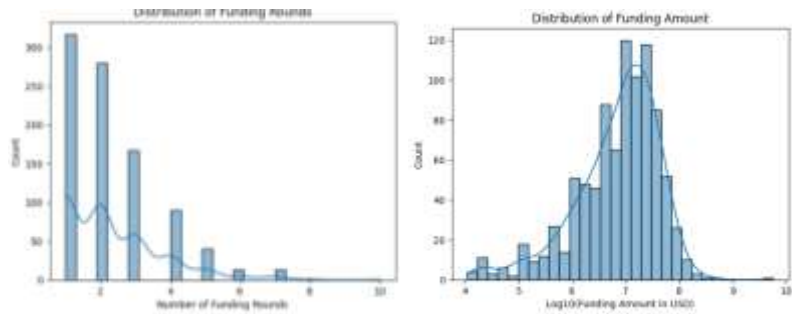| | |
|---|---|
| Date | 18 June 2025 |
| Team ID | SWTID1749880888 |
| Project Title | Prosperity Prognosticator: Machine Learning for Startup Success Prediction |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

| Section | Description |
|---|---|
| Data Overview | Dimension: 923 rows x 49 columns<br>Descriptive Statistics:<br> |
| Univariate Analysis |  |

Bivariate Analysis

| | |
|---|---|
| Multivariate Analysis |  |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data |  |
| Dropping Irrelevant Columns |  |

| | |
|---|---|
| Feature Engineering |  |
| Save Processed Data |  |