# Behavioral Data Analysis Report



**1. Introduction**

This report aims to analyze behavioral data to understand patterns and factors affecting the price of electric vehicles (EVs). The analysis involves data preprocessing, exploratory data analysis (EDA), feature engineering, and predictive modeling using machine learning techniques.

**2. Data Preprocessing**

**2.1 Data Cleaning**

- **Handling Missing Values:** The dataset is checked for missing values, and appropriate measures are taken to handle them, such as imputation or removal.

- **Categorical Encoding:** Categorical features are encoded using techniques like OneHotEncoder to convert them into numerical format.

**2.2 Feature Scaling**

- Numerical features are scaled using StandardScaler to ensure they have a mean of 0 and a standard deviation of 1. This step is crucial for models that rely on distance metrics.

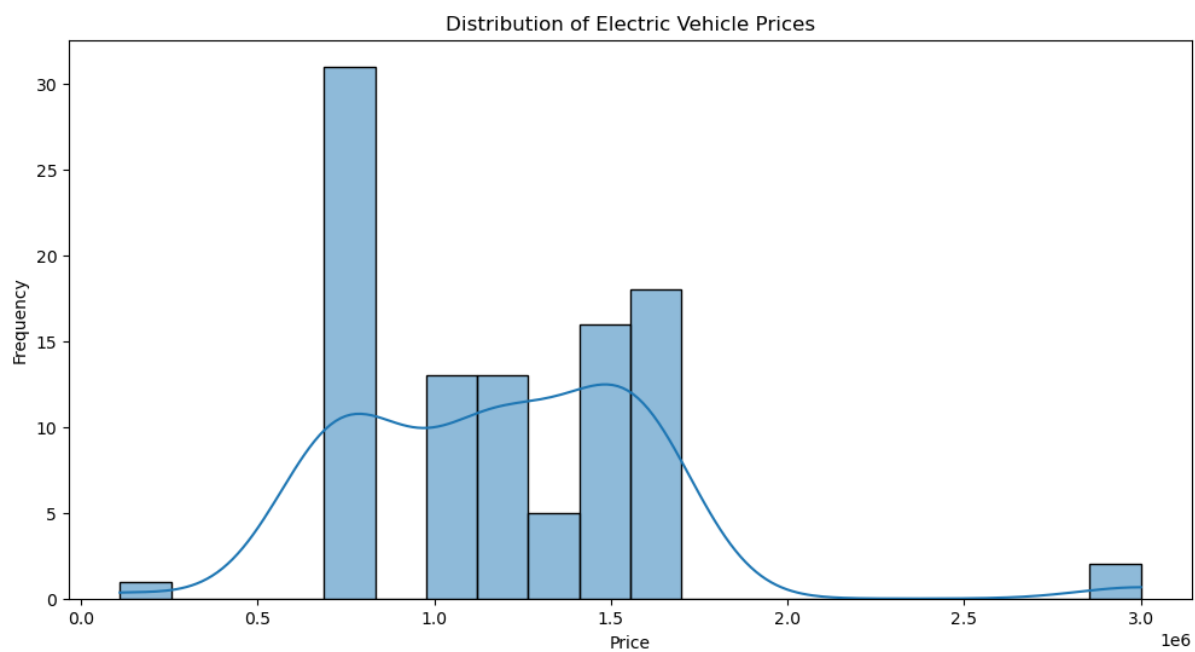**3. Exploratory Data Analysis (EDA)**

**3.1 Descriptive Statistics**

- The dataset's statistical summary is provided, including measures like mean, median, standard deviation, and interquartile ranges for numerical features.
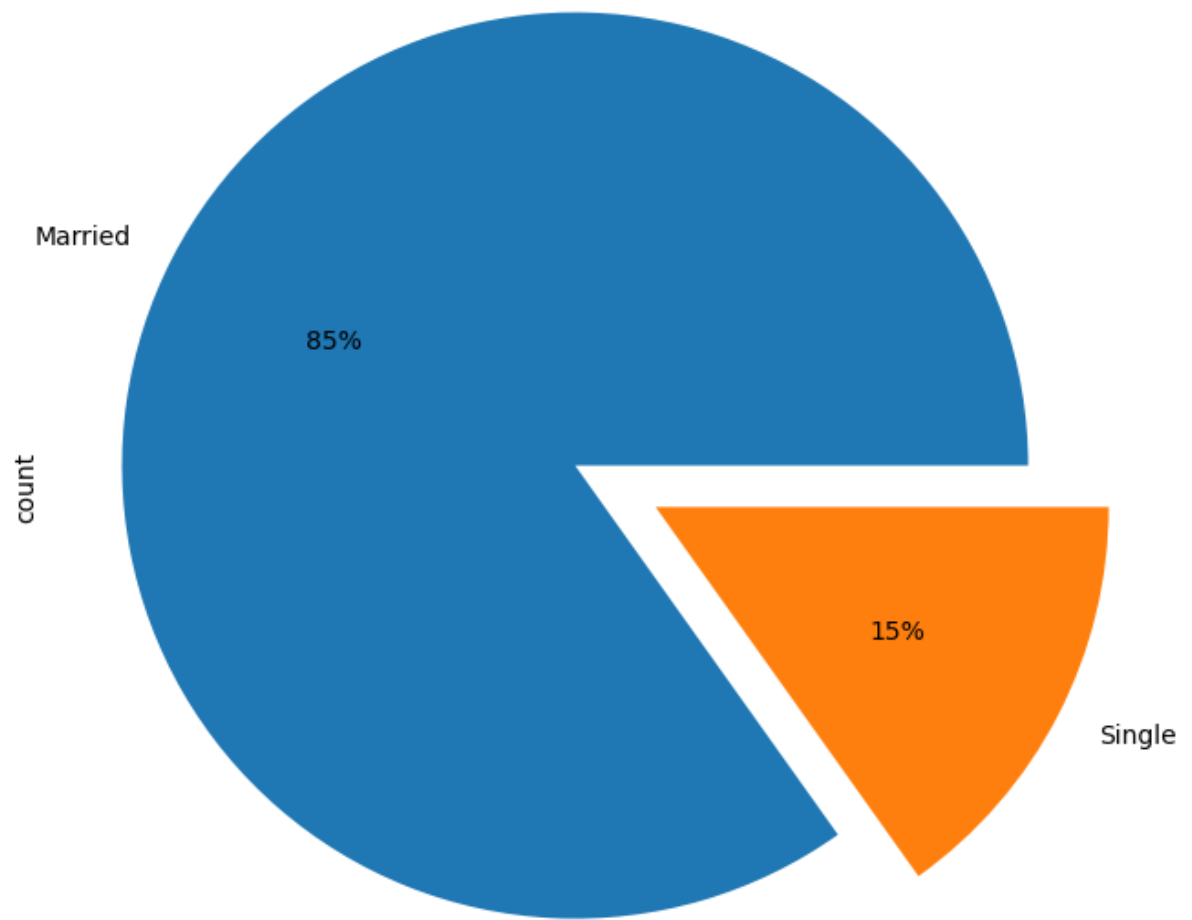
**3.2 Distribution Analysis**

- Histograms and boxplots are used to visualize the distribution of key numerical features like age, total salary, and the target variable, Price.
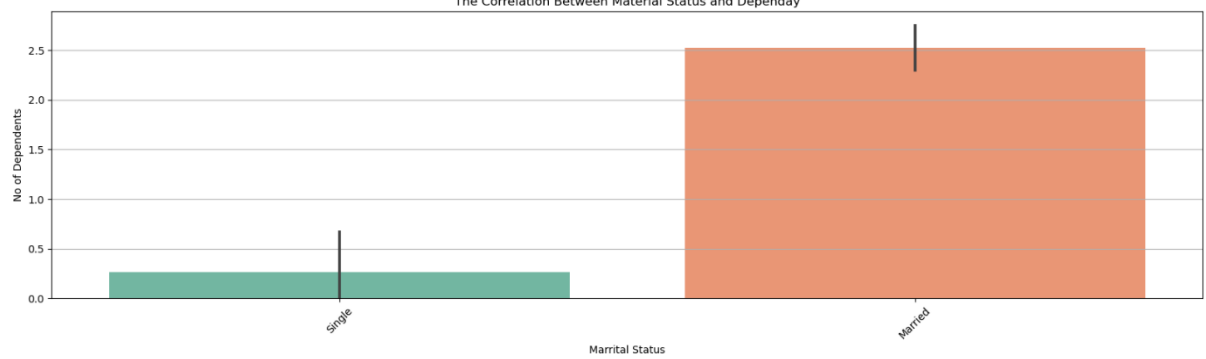
- The distribution of categorical features, such as profession, marital status, education, and personal loan status, is also examined.
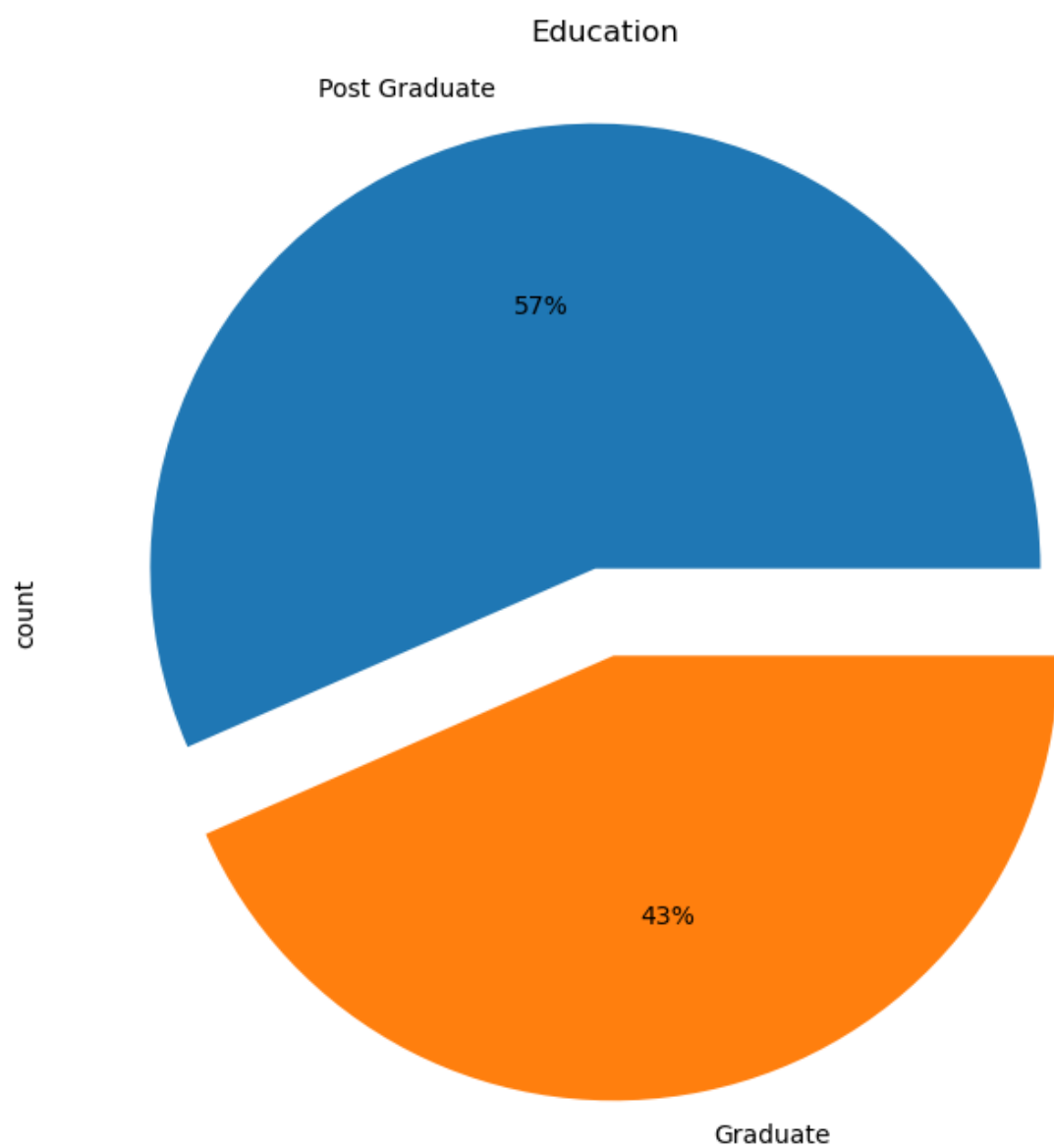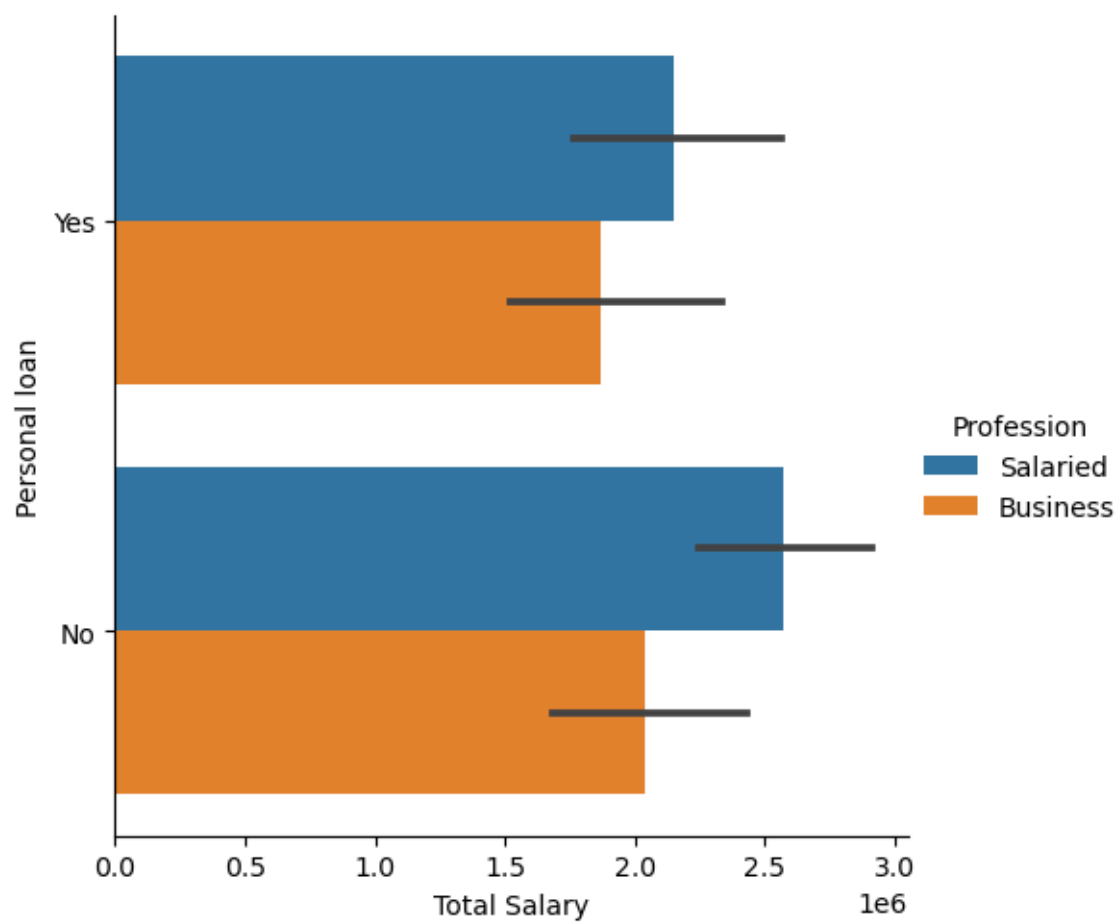
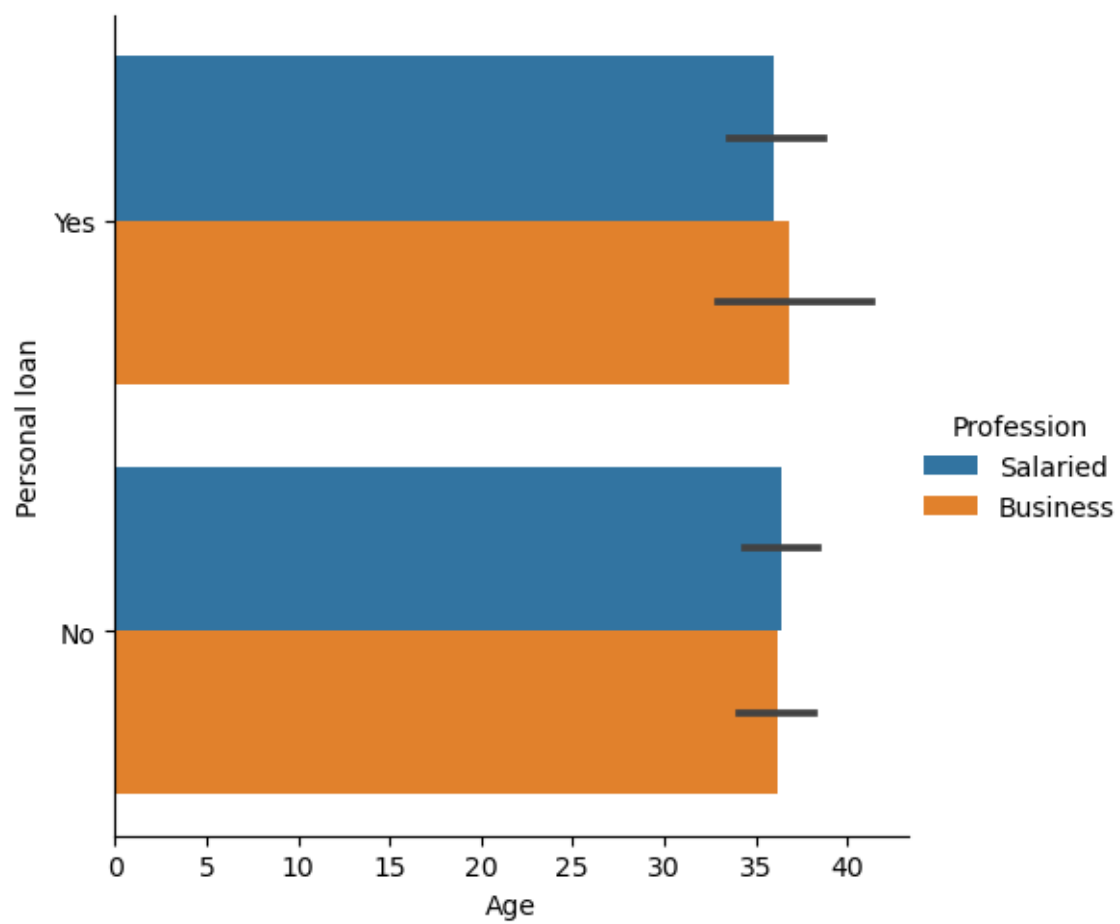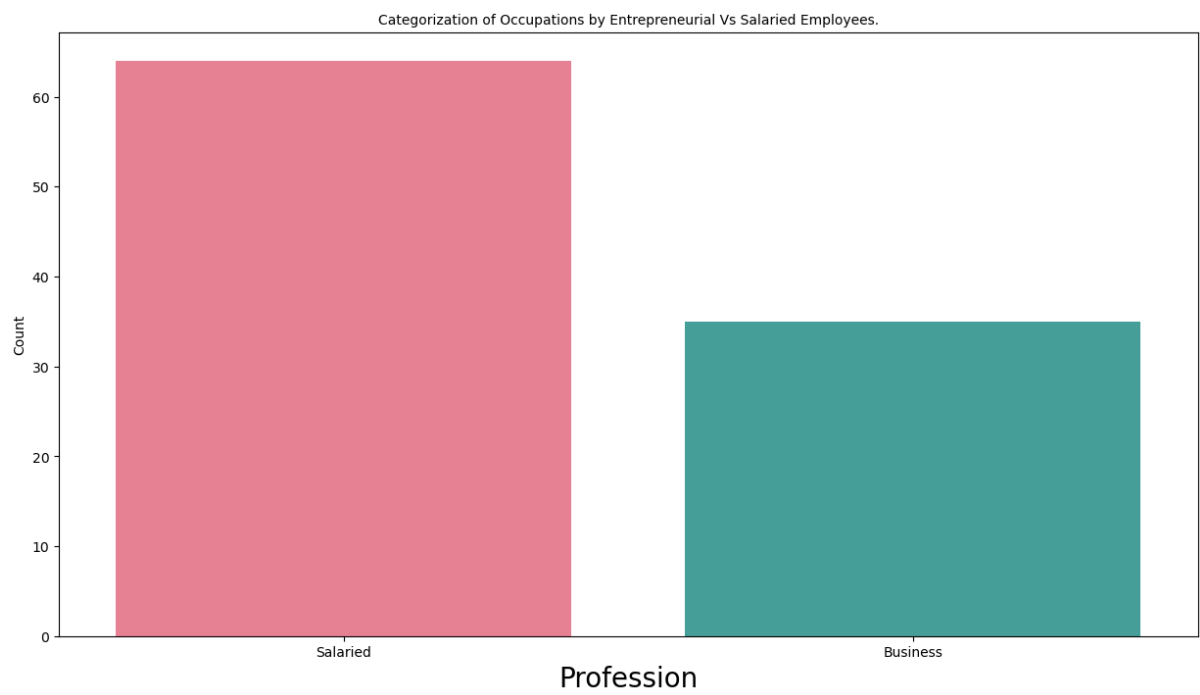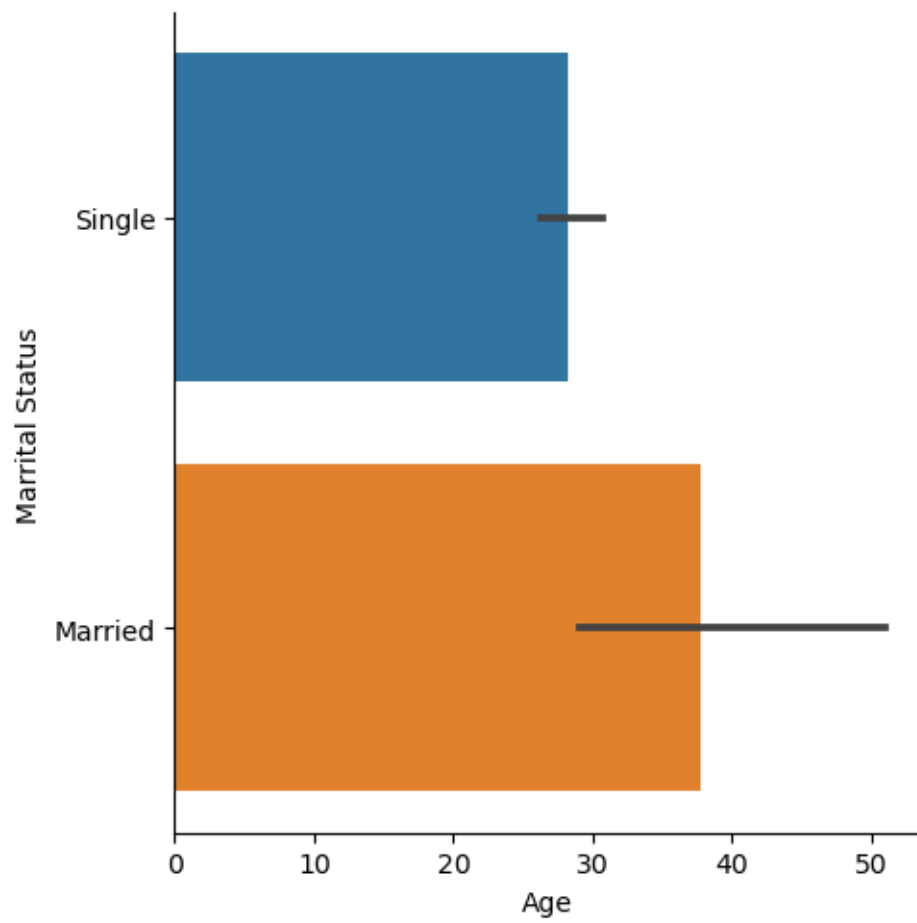Distribution of Electric Vehicle Prices

# Marrital Status



The Correlation Between Material Status and Dependay

# Education
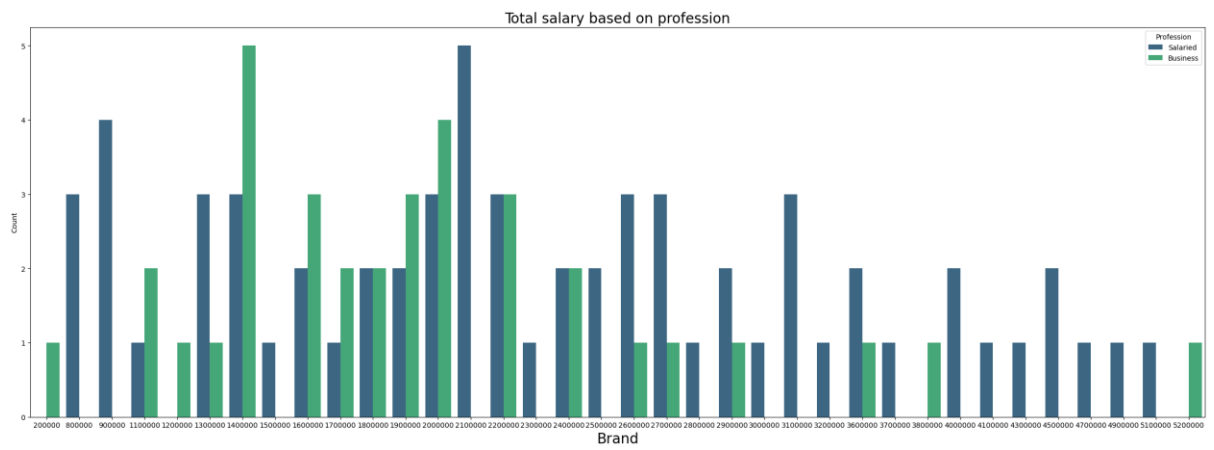
Post Graduate

57%

count

43%

Graduate

Categorization of Occupations by Entrepreneurial Vs Salaried Employees.

Profession with Age



EV Based on Brand



Total salary based on profession
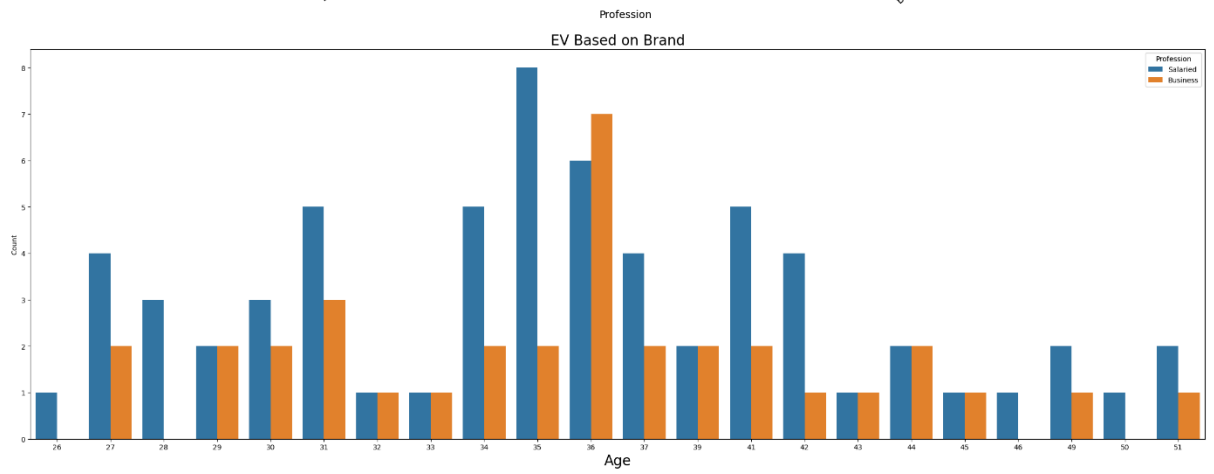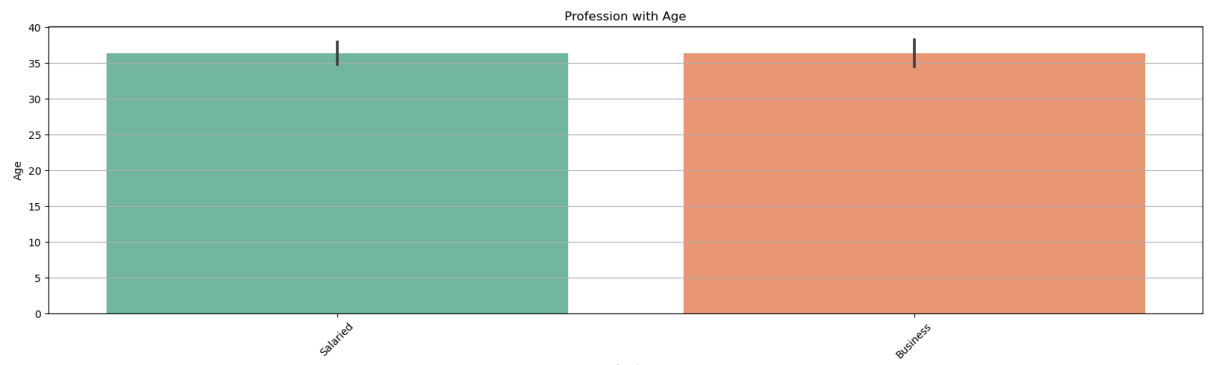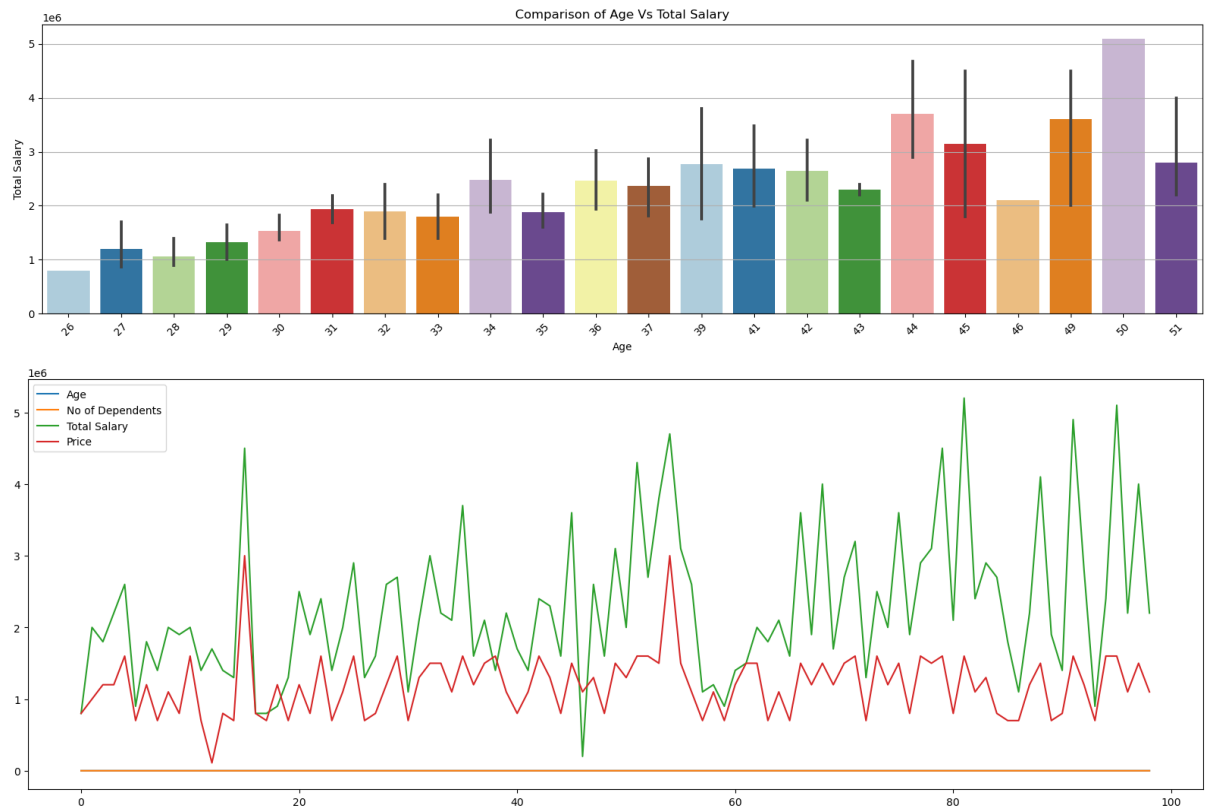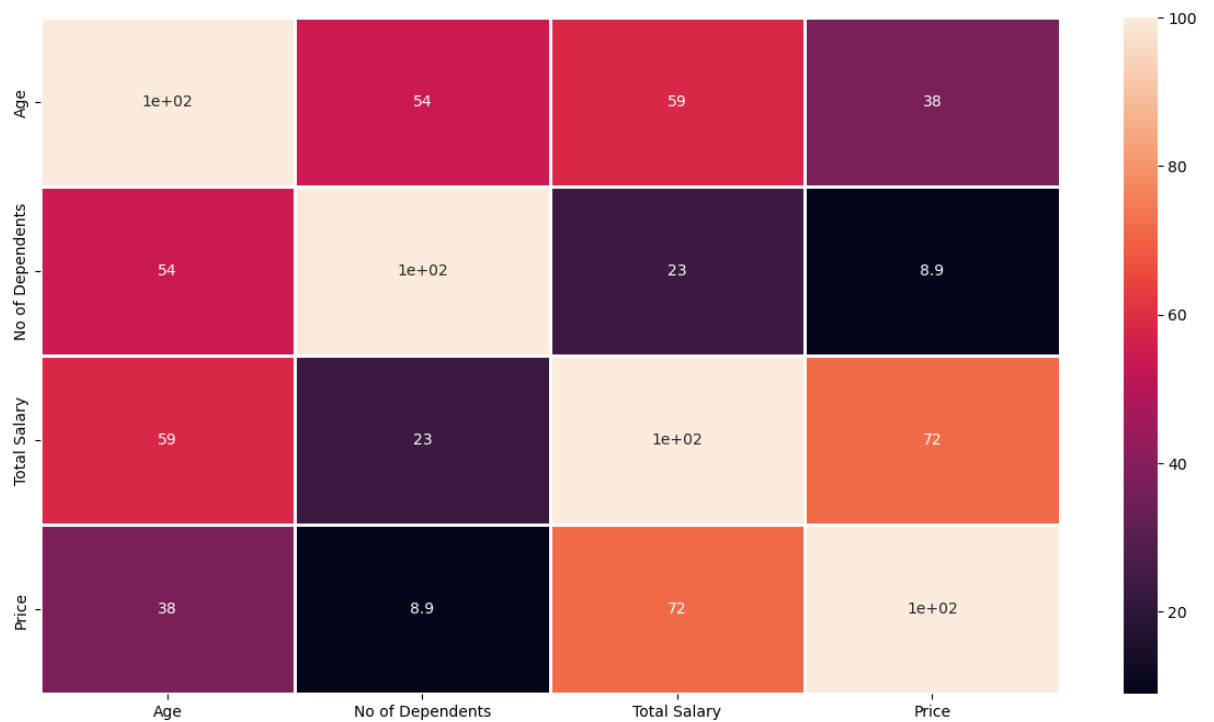
Comparison of Age Vs Total Salary

### 3.3 Correlation Analysis

- The Pearson correlation coefficient is calculated to identify relationships between numerical features.

- Heatmaps are used to visualize the correlation matrix and highlight significant correlations.



### 3.4 Feature Importance

- Feature importance is assessed using models like Random Forest to understand the impact of each feature on the target variable.

## 4. Predictive Modeling

### 4.1 Model Selection

- **Random Forest Regressor** is chosen for its robustness and ability to handle both numerical and categorical features.

### 4.2 Model Training

- The data is split into training and testing sets.

- A pipeline is created to streamline preprocessing and model training.

### 4.3 Hyperparameter Tuning

- GridSearchCV is used for hyperparameter tuning to find the best combination of parameters for the Random Forest model.

### 4.4 Model Evaluation

- The model's performance is evaluated using metrics like Mean Squared Error (MSE) and R-squared ($R^2$) score.

**Results:**

- **Random Forest Mean Squared Error:** 67376721049.65011

- **Random Forest R² Score:** 0.3493664276321341

## 5. Inferences and Recommendations

### 5.1 Key Findings

- **Feature Impact:** Total Salary and Age are significant predictors of EV Price.

- **Model Performance:** The Random Forest model achieved a reasonable $R^2$ score, indicating a moderate fit to the data.

### 5.2 Recommendations

- **Data Quality:** Improve data quality by addressing missing values and potential outliers more rigorously.

- **Feature Engineering:** Explore additional features and interaction terms to enhance model performance.

- **Advanced Models:** Consider using more advanced models or ensemble techniques to improve prediction accuracy.

## 6. Conclusion

This analysis provided valuable insights into the factors influencing EV prices. While the current model offers a moderate level of predictive accuracy, further refinement and exploration of additional features could enhance its performance.