# A Comprehensive Guide to SEMMA Methodology: Case Study on the 100 Largest Companies

Venkatesh Makkena

# Contents

**Abstract**

This research paper aims to provide a thorough application of the SEMMA (Sample, Explore, Modify, Model, Assess) methodology for data science on a dataset of the 100 largest companies. This in-depth guide includes detailed theoretical discussions, practical code snippets, and their respective outputs.

# 1    Introduction

Data science has become a cornerstone in various sectors, enabling businesses to derive valuable insights and make informed decisions. However, the transformation from raw data to actionable intelligence is often a complex process, requiring a structured approach. The SEMMA methodology offers a reliable framework for navigating this intricate landscape. This paper serves as a rigorous application of SEMMA to a real-world dataset of the 100 largest companies.

# 2    Methods

The SEMMA methodology, standing for Sample, Explore, Modify, Model, and Assess, is a structured approach to data science. Each step plays a pivotal role in the data science pipeline, ensuring the effectiveness and robustness of the final model.

## 2.1    Sample

Sampling is often the foundational step in data science. The primary goal is to take a manageable subset of the dataset to expedite initial computations. By working with a subset, one can perform quick, exploratory data analyses to gauge the data's structure, missing values, and potential outliers.

```
# Load the dataset
df = pd.read_csv('/mnt/data/Largest_Companies.csv')
# Display first few rows
df.head()
```

Output:

```
      Company  Rank  Revenue (USD millions)  Revenue growth  Employees
0     Walmart     1                523964.0             1.9    2200000
1      Amazon     2                280522.0            20.5     798000
2    McKesson     3                214319.0             2.6      70000
3       Cigna     4                153566.0           209.4      73700
```

## 2.2 Explore

The Explore phase aims to provide a deeper understanding of the dataset. This involves generating summary statistics and identifying missing values. It also includes the vital step of data visualization, which can help identify trends, patterns, and outliers in the data.

```
# Summary statistics
df.describe()
# Missing values
df.isnull().sum()
```

Output:

```
Summary Statistics:
        Rank  Revenue (USD millions)  Revenue growth      Employees
count  100.0              100.000000      100.000000   1.000000e+02
mean    50.5            69667.820000       13.610000   1.230200e+05
std     29.0            97797.902327       29.487557   3.522456e+05
min      1.0            24783.000000      -62.300000   3.100000e+01
max    100.0           523964.000000      209.400000   2.200000e+06

Missing Values:
Company                   0
Rank                      0
Revenue (USD millions)    0
Revenue growth            0
Employees                 0
```

## 2.3 Modify

The Modify phase plays a significant role in enhancing the model's performance. It involves data cleaning, preprocessing, and feature selection. Fea-

ture selection is especially crucial as it helps the model focus on the most relevant variables, improving both performance and interpretability.

```
# Remove outliers
df_clean = df[(z_scores < 3).all(axis=1)]
# Data Preprocessing
df_clean_encoded = pd.get_dummies(df_clean, columns=categorical_cols)
```

Output:

```
Rows Before: 100
Rows After Outlier Removal: 96
```

## 2.4 Model

The Model phase is the core of the data science pipeline. It involves the selection and training of machine learning algorithms to build predictive models. The importance of starting with a baseline model cannot be overstated as it provides a benchmark for evaluating more complex models.

```
# Baseline model
baseline_model = LinearRegression()
baseline_model.fit(X_train, y_train)
# Advanced models
rf_model = RandomForestRegressor()
rf_model.fit(X_train, y_train)
```

Output:

```
Baseline Model R2: 0.483
Random Forest Model R2: 0.968
```

## 2.5 Assess

The Assess phase involves evaluating the model's performance based on specific metrics and making final recommendations. This phase is critical as it dictates the model's deployment and its utility in decision-making processes.

```
# Model Comparison
model_comparison = pd.DataFrame({
    'Model': ['Baseline', 'Random-Forest', 'Gradient-Boosting'],
```

4

```
    'R2': [0.483, 0.968, 0.972]
})
```

Output:

```
Model Comparison:
              Model     R2
0          Baseline  0.483
1     Random Forest  0.968
2 Gradient Boosting  0.972
```

# 3  Results

The Random Forest and Gradient Boosting models showed a remarkable improvement over the baseline model. In particular, the Gradient Boosting model was highly effective, explaining approximately 97.2% of the variance in the dataset.

# 4  Conclusion

This research paper serves as a detailed, comprehensive guide to applying the SEMMA methodology to a real-world dataset of the 100 largest companies. From the importance of initial data sampling to the intricacies of model assessment, this study has demonstrated the robustness and effectiveness of the SEMMA methodology.

# 5  References

Dataset sourced from Kaggle: `https://www.kaggle.com/datasets/omikumarmakadia2121/100-largest-companies`