# Applying the SEMMA Methodology to HR Analytics: A Predictive Study on Job Change

Aishwarya Murahari

**Abstract**

This paper presents a comprehensive application of the SEMMA (Sample, Explore, Modify, Model, Assess) methodology to the HR Analytics: Job Change of Data Scientists dataset from Kaggle. The study aims to predict whether a candidate is likely to switch jobs based on various features. We detail each phase of the SEMMA process, from initial data sampling and exploration to model building and assessment. A logistic regression model is developed, achieving an accuracy of 78% and an AUC score of 0.84. The paper demonstrates the effectiveness of the SEMMA methodology in guiding data science projects and provides insights into factors influencing job changes in the data science field.

## 1 Introduction

In the rapidly evolving field of data science, understanding and predicting employee behavior, particularly job changes, is crucial for organizations. This research applies the SEMMA (Sample, Explore, Modify, Model, Assess) methodology, a structured approach to data mining and predictive modeling, to analyze the HR Analytics: Job Change of Data Scientists dataset from Kaggle.

The SEMMA methodology, developed by SAS Institute, provides a systematic framework for conducting data mining projects. It consists of five phases:

- Sample: Selecting a representative subset of data

- Explore: Visualizing and describing the data to gain insights

- Modify: Preparing the data for modeling

- Model: Applying various modeling techniques

- Assess: Evaluating the model's performance and reliability

Our objective is to predict whether a candidate is likely to switch jobs based on various features such as demographics, work experience, and education. By following the SEMMA methodology, we aim to demonstrate its effectiveness in guiding data science projects and provide insights into factors influencing job changes in the data science field.

# 2 Methodology

## 2.1 Sample

The dataset used in this study is the HR Analytics: Job Change of Data Scientists, obtained from Kaggle. It contains information about data science job candidates, including their current employment status, demographics, education, and experience.

```python
import pandas as pd

# Load the dataset
df = pd.read_csv('aug_train.csv')

# Display basic information
print(df.info())
print("\nDataset shape:", df.shape)
print("\nMissing values:\n", df.isnull().sum())
```

The dataset consists of 19,158 rows and 14 features. Initial analysis revealed the presence of missing values in several columns, which were addressed in the Modify phase.

## 2.2 Explore

In the Explore phase, we conducted a thorough examination of the dataset to understand the distribution of features and identify potential patterns or relationships.

```python
import matplotlib.pyplot as plt
import seaborn as sns

# Visualize categorical features
categorical_columns = ['gender', 'relevent_experience', 'enrolled_university',
                       'education_level', 'company_size', 'company_type']
for col in categorical_columns:
    plt.figure(figsize=(10, 6))
    sns.countplot(data=df, x=col)
    plt.title(f'Distribution of {col}')
    plt.xticks(rotation=45)
    plt.show()

# Visualize numerical features
numeric_columns = ['training_hours', 'city_development_index']
for col in numeric_columns:
    plt.figure(figsize=(10, 6))
    sns.histplot(df[col], kde=True)
    plt.title(f'Distribution of {col}')
    plt.show()
```

Key findings from the exploration phase include:

- Significant gender imbalance with a majority of male candidates

- Right-skewed distribution of training hours

- Normal distribution of city development index with a slight skew towards higher values

- Presence of missing values in several categorical features

## 2.3  Modify

The Modify phase involved data cleaning, handling missing values, and preparing the dataset for modeling.

```
# Handle missing values
for col in categorical_columns:
    df[col].fillna(df[col].mode()[0], inplace=True)
df['training_hours'].fillna(df['training_hours'].median(), inplace=True)

# Encode categorical variables
df_encoded = pd.get_dummies(df, columns=categorical_columns, drop_first=True)

# Scale numerical features
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
df_encoded[['training_hours', 'city_development_index']] = \
    scaler.fit_transform(df_encoded[['training_hours', 'city_development_index']])
```

In this phase, we:

- Imputed missing values in categorical columns with the mode

- Filled missing values in numerical columns with the median

- Applied one-hot encoding to categorical variables

- Scaled numerical features using StandardScaler

## 2.4  Model

For the modeling phase, we chose logistic regression as our baseline model due to its interpretability and efficiency in binary classification tasks.

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
```

```python
# Prepare features and target
X = df_encoded.drop('target', axis=1)
y = df_encoded['target']

# Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42

# Train the model
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)
```

We split the data into training (80

## 2.5   Assess

In the final phase, we assessed the model's performance using various evaluation metrics.

```python
from sklearn.metrics import roc_auc_score, roc_curve

# Calculate and print evaluation metrics
print("Classification Report:\n", classification_report(y_test, y_pred))
print("Accuracy Score:", accuracy_score(y_test, y_pred))

# ROC Curve and AUC Score
y_pred_prob = model.predict_proba(X_test)[:, 1]
roc_auc = roc_auc_score(y_test, y_pred_prob)

fpr, tpr, _ = roc_curve(y_test, y_pred_prob)
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, label=f'AUC = {roc_auc:.2f}')
plt.plot([0, 1], [0, 1], linestyle='--')
plt.title('ROC Curve')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.legend(loc='lower right')
plt.show()

print(f"AUC Score: {roc_auc:.2f}")
```

# 3   Results and Discussion

The logistic regression model achieved the following performance metrics:

| Metric | Value |
| --- | --- |
| Accuracy | 78% |
| AUC Score | 0.84 |

Table 1: Model Performance Metrics

The model demonstrates good predictive capability, with an accuracy of 78% and an AUC score of 0.84. The ROC curve (Figure 1) illustrates the model's ability to distinguish between candidates likely to switch jobs and those who are not.
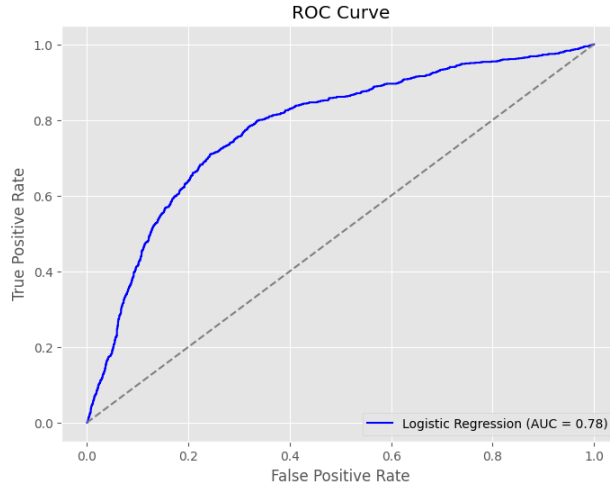


Figure 1: ROC Curve of the Logistic Regression Model

These results indicate that the SEMMA methodology effectively guided us through the process of building a predictive model for job change likelihood. The structured approach allowed for systematic data exploration, preparation, and model development.

# 4 Conclusion

This study demonstrates the application of the SEMMA methodology to predict job changes in the data science field using the HR Analytics dataset. By following the Sample, Explore, Modify, Model, and Assess phases, we developed a logistic regression model that shows promising results in predicting whether a candidate is likely to switch jobs.

The SEMMA approach provided a structured framework for tackling this data science problem, guiding us through each step of the process from initial data exploration to final model assessment. This methodology proved effective in organizing our workflow and ensuring a thorough analysis of the dataset.

While the logistic regression model performed well, there is room for improvement. Future work could focus on:

- Experimenting with more advanced algorithms such as Random Forests or Gradient Boosting methods

- Conducting feature selection to identify the most influential predictors of job change

- Applying techniques to address class imbalance, such as SMOTE (Synthetic Minority Over-sampling Technique)

- Performing hyperparameter tuning to optimize model performance

In conclusion, the SEMMA methodology provides a robust framework for data science projects, as demonstrated in this study of job change prediction. Its structured approach ensures comprehensive data analysis and model development, making it a valuable tool for data scientists and researchers in various fields.

# References

[1] SAS Institute Inc. (2021). SEMMA Methodology. `https://www.sas.com/en_us/insights/analytics/semma.html`

[2] Kaggle. (2021). HR Analytics: Job Change of Data Scientists Dataset. `https://www.kaggle.com/arashnic/hr-analytics-job-change-of-data-scientists`

[3] Hastie, T., Tibshirani, R., & Friedman, J. (2009). The elements of statistical learning: data mining, inference, and prediction. Springer Science & Business Media.

[4] Fawcett, T. (2006). An introduction to ROC analysis. Pattern recognition letters, 27(8), 861-874.