# Sheth L.U.J College of Arts & Sir M.V. College of Science and Commerce
## Data Science
## PRACTICAL NO. 6

**Aim: Regression and Its Types**

   a) Implement simple linear regression using a dataset.
   b) Explore and interpret the regression model coefficients and goodness-of-fit measures.
   c) Extend the analysis to multiple linear regression and assess the impact of additional predictors.

## CODE:

> ### *Importing libraries*

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Lasso, Ridge
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.metrics import mean_squared_error, r2_score
```

> ### *Load Dataset*

```
df = pd.read_csv('insurance.csv')
print("Dataset Overview:\n")
display(df.head())
```

> ### *Summary statistics*

```
print("\nSummary Statistics:\n")
display(df.describe())
```

> ### *Check data types*

```
print("\nData Types:\n")
display(df.dtypes)
```

> ### *Data Visualization*

```
# Distribution of target variable (charges)
plt.figure(figsize=(8,5))
sns.histplot(df['charges'], bins=20, kde=True, color='skyblue')
plt.title('Distribution of Charges')
plt.xlabel('Charges')
plt.ylabel('Frequency')
plt.show()

# Relationship between age and charges
plt.figure(figsize=(8,5))
sns.scatterplot(x='age', y='charges', data=df, hue='smoker', style='sex', s=100)
plt.title('Age vs Charges')
plt.show()

# Pairplot for numerical variables
sns.pairplot(df, hue='smoker')
plt.show()
```

Abin Pillai T104                                                                                    Data Science

```python
# Correlation heatmap
plt.figure(figsize=(8,6))
numeric_df = df.select_dtypes(include=np.number)
sns.heatmap(numeric_df.corr(), annot=True, cmap='coolwarm')
plt.title('Correlation Heatmap (Numeric Features)')
plt.show()
```

### ➢ *Simple Linear Regression (SLR)*

```python
# Split dataset into independent (X) and dependent (y) variables
X = df[['age']].values  # Age independent variable
y = df['charges'].values  # Charges dependent variable

# Split into train/test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train linear regression model
slr = LinearRegression()
slr.fit(X_train, y_train)

# Predict
y_train_pred = slr.predict(X_train)
y_test_pred = slr.predict(X_test)

# Model coefficients
print("Simple Linear Regression Coefficients:")
print(f"Slope (m): {slr.coef_[0]:.2f}")
print(f"Intercept (c): {slr.intercept_:.2f}")

# Plot training set vs predictions
plt.figure(figsize=(8,5))
plt.scatter(X_train, y_train, color='blue', label='Training Data')
plt.plot(X_train, y_train_pred, color='red', linewidth=2, label='Regression Line')
plt.title('SLR: Training Set - Age vs Charges')
plt.xlabel('Age')
plt.ylabel('Charges')
plt.legend()
plt.show()

# Plot test set vs predictions
plt.figure(figsize=(8,5))
plt.scatter(X_test, y_test, color='green', label='Test Data')
plt.plot(X_test, y_test_pred, color='red', linewidth=2, label='Predictions')
plt.title('SLR: Test Set - Age vs Charges')
plt.xlabel('Age')
plt.ylabel('Charges')
plt.legend()
plt.show()

# Performance metrics
print("SLR Performance on Test Set:")
print(f"R-squared: {r2_score(y_test, y_test_pred):.2f}")
print(f"Mean Squared Error: {mean_squared_error(y_test, y_test_pred):.2f}")
```

### ➢ *Multiple Linear Regression (MLR)*

```python
# Independent variables: all except charges
```

Abin Pillai T104                                                                 Data Science

```python
X = df.iloc[:, :-1].values
y = df['charges'].values

# One-hot encoding for categorical variables: 'sex', 'smoker', 'region'
ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [1,4,5])], remainder='passthrough')
X = np.array(ct.fit_transform(X))

# Split into train/test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train Multiple Linear Regression model
mlr = LinearRegression()
mlr.fit(X_train, y_train)

# Predict
y_pred = mlr.predict(X_test)

# Performance metrics
print("Multiple Linear Regression Performance:")
print(f"R-squared: {r2_score(y_test, y_pred):.2f}")
print(f"Mean Squared Error: {mean_squared_error(y_test, y_pred):.2f}")

# Plot Actual vs Predicted
plt.figure(figsize=(8,5))
plt.scatter(y_test, y_pred, color='purple')
plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red', linestyle='--')
plt.title('MLR: Actual vs Predicted Charges')
plt.xlabel('Actual Charges')
plt.ylabel('Predicted Charges')
plt.show()

# Lasso Regression (Regularization)
lasso = Lasso(alpha=0.1)
lasso.fit(X_train, y_train)
y_lasso_pred = lasso.predict(X_test)
print("Lasso Regression R-squared:", r2_score(y_test, y_lasso_pred))

# Ridge Regression (Regularization)
ridge = Ridge(alpha=1.0)
ridge.fit(X_train, y_train)
y_ridge_pred = ridge.predict(X_test)
print("Ridge Regression R-squared:", r2_score(y_test, y_ridge_pred))

# Elastic Net Regression
from sklearn.linear_model import ElasticNet
elastic = ElasticNet(alpha=0.1, l1_ratio=0.5)  # l1_ratio=0.5 balances L1 and L2
elastic.fit(X_train, y_train)
y_elastic_pred = elastic.predict(X_test)
print("Elastic Net R-squared:", r2_score(y_test, y_elastic_pred))
```

## Output:

# Sheth L.U.J College of Arts & Sir M.V. College of Science and Commerce
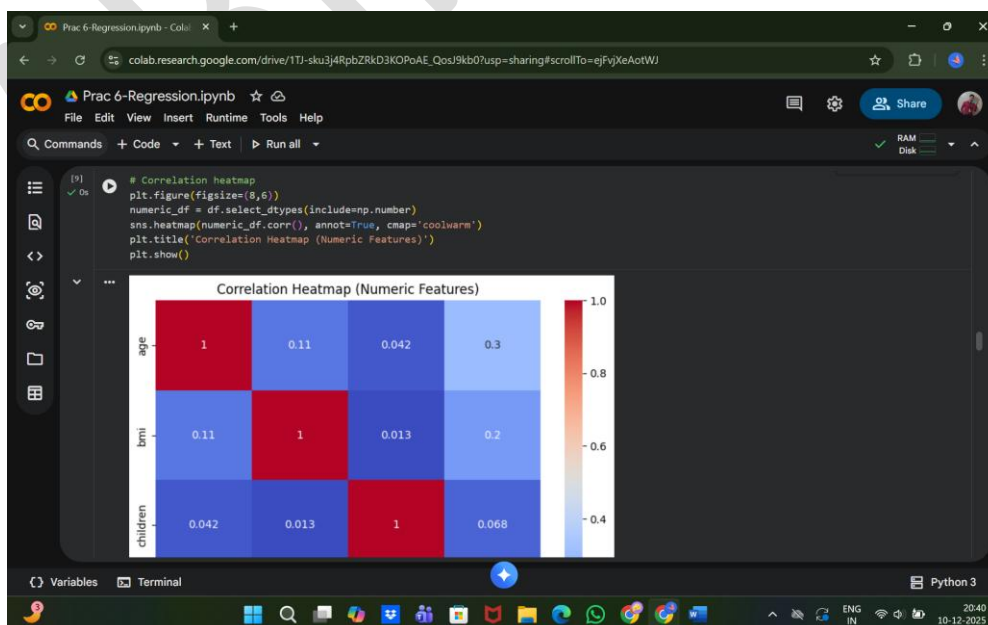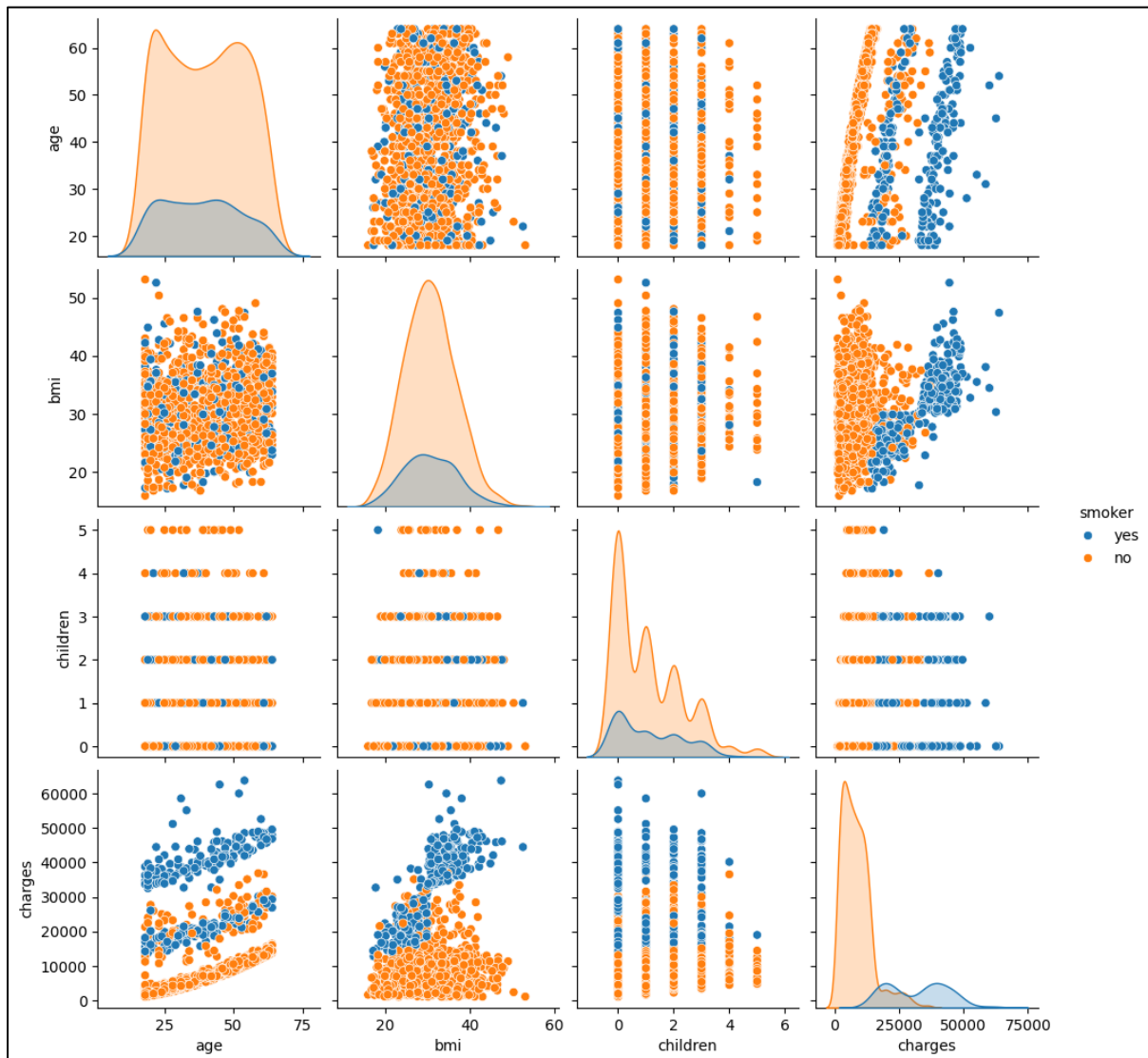## Data Science
## PRACTICAL NO. 6

# Sheth L.U.J College of Arts & Sir M.V. College of Science and Commerce
## Data Science
## PRACTICAL NO. 6





Distribution of Charges

# Sheth L.U.J College of Arts & Sir M.V. College of Science and Commerce
## Data Science
## PRACTICAL NO. 6

# Sheth L.U.J College of Arts & Sir M.V. College of Science and Commerce
## Data Science
## PRACTICAL NO. 6



Correlation Heatmap (Numeric Features)

The test set plot confirms this, showing that many predicted values (red line) underestimate or overestimate actual charges, especially for higher values. Together, these plots highlight the limitation of simple linear regression with a single predictor and the need for additional variables in a multiple regression model to improve prediction accuracy.

```python
# Performance metrics
print("SLR Performance on Test Set:")
print(f"R-squared: {r2_score(y_test, y_test_pred):.2f}")
print(f"Mean Squared Error: {mean_squared_error(y_test, y_test_pred):.2f}")
```

```
SLR Performance on Test Set:
R-squared: 0.12
Mean Squared Error: 135983957.48
```

The SLR model using only age as a predictor gives an R-squared of 0.12, which means that only 12% of the variation in insurance charges is explained by age alone. The high Mean Squared Error (≈136 million) indicates that the predictions are not very accurate, especially for individuals with higher charges. This shows that age alone is not sufficient to predict insurance costs, and other factors like BMI, smoker status, children, and region must be included to improve the model.

## Multiple Linear Regression (MLR)



```python
# Independent variables: all except charges
X = df.iloc[:, :-1].values
y = df['charges'].values
```

```python
# One-hot encoding for categorical variables: 'sex', 'smoker', 'region'
ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [1,4,5])], remainder='passthrough')
X = np.array(ct.fit_transform(X))
```

```python
# Split into train/test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```python
# Train Multiple Linear Regression model
mlr = LinearRegression()
mlr.fit(X_train, y_train)
```
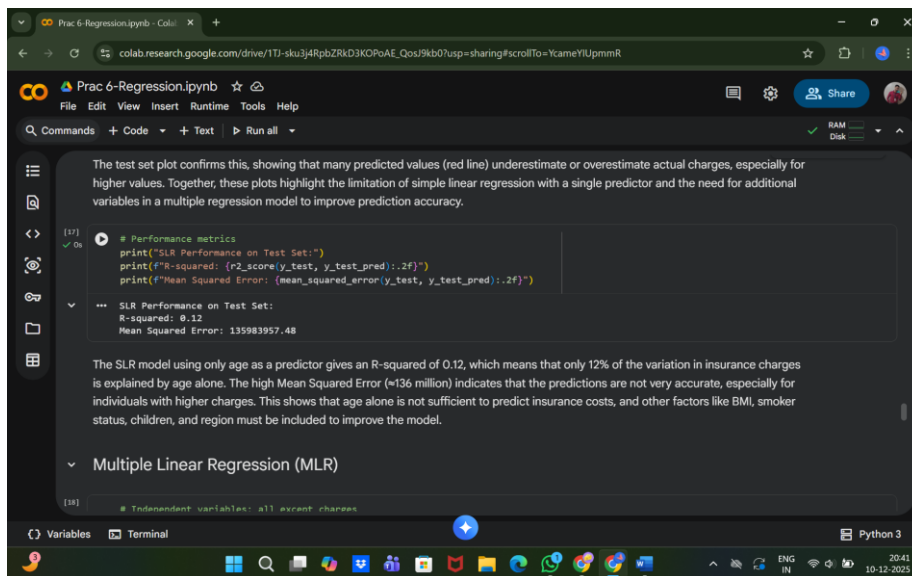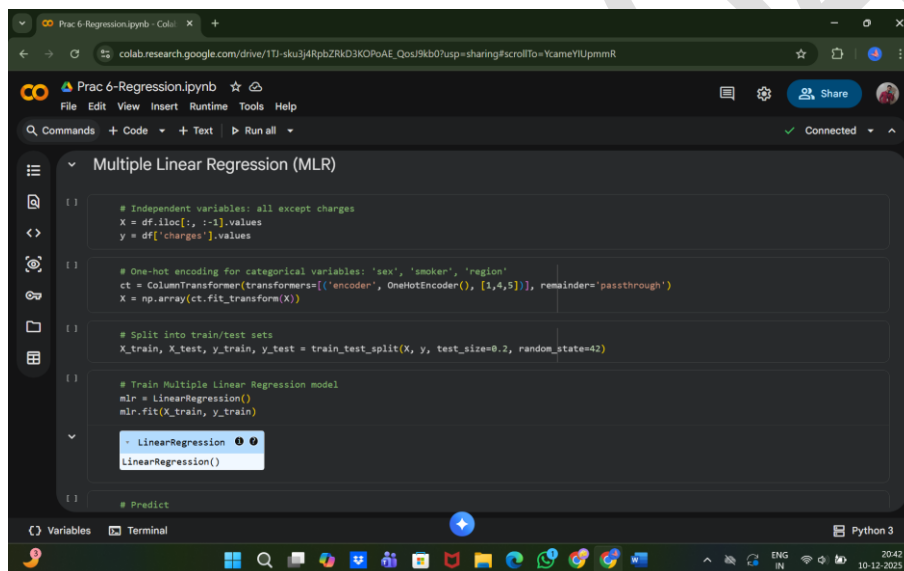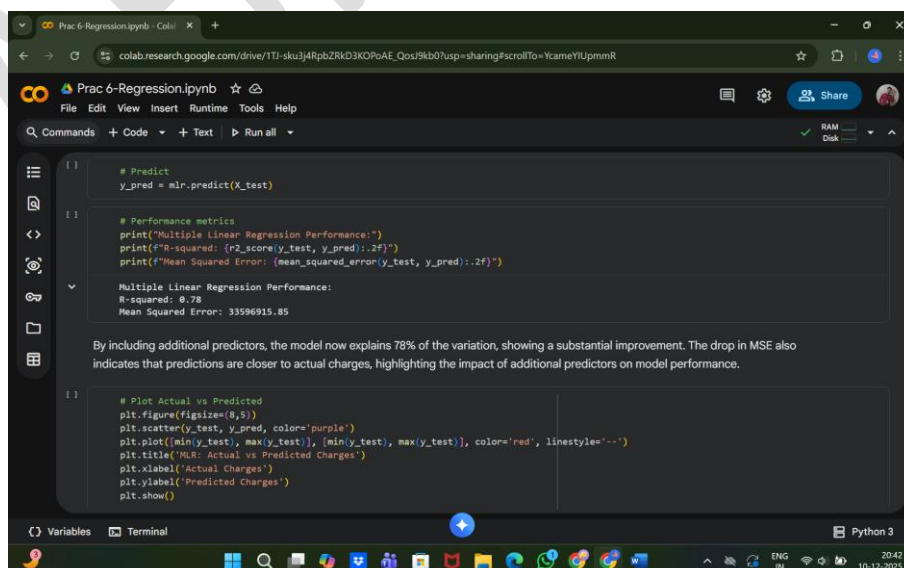
```
LinearRegression
LinearRegression()
```



```python
# Predict
y_pred = mlr.predict(X_test)
```

```python
# Performance metrics
print("Multiple Linear Regression Performance:")
print(f"R-squared: {r2_score(y_test, y_pred):.2f}")
print(f"Mean Squared Error: {mean_squared_error(y_test, y_pred):.2f}")
```

```
Multiple Linear Regression Performance:
R-squared: 0.78
Mean Squared Error: 33596915.85
```
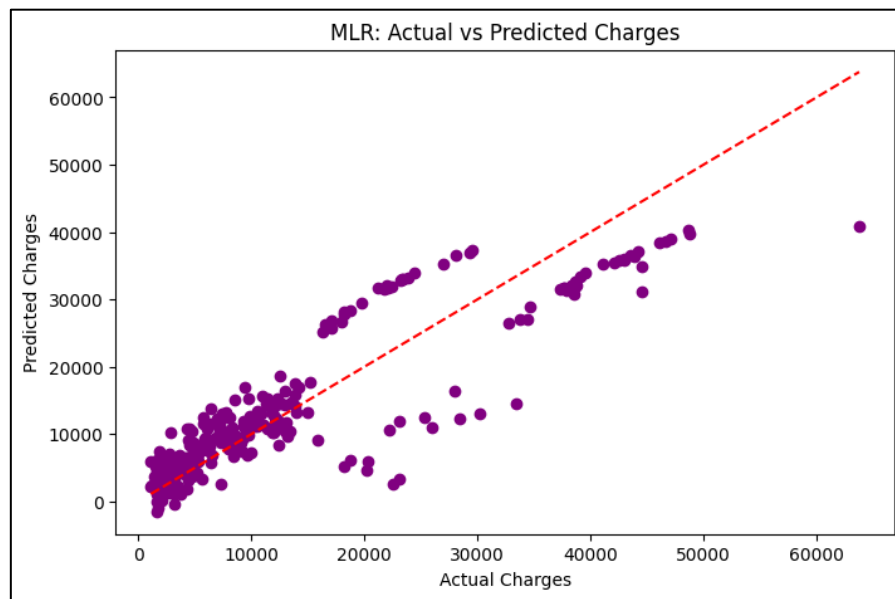
By including additional predictors, the model now explains 78% of the variation, showing a substantial improvement. The drop in MSE also indicates that predictions are closer to actual charges, highlighting the impact of additional predictors on model performance.

```python
# Plot Actual vs Predicted
plt.figure(figsize=(8,5))
plt.scatter(y_test, y_pred, color='purple')
plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red', linestyle='--')
plt.title('MLR: Actual vs Predicted Charges')
plt.xlabel('Actual Charges')
plt.ylabel('Predicted Charges')
plt.show()
```

# Sheth L.U.J College of Arts & Sir M.V. College of Science and Commerce
## Data Science
## PRACTICAL NO. 6



MLR: Actual vs Predicted Charges



The scatter plot shows that the multiple linear regression model predicts lower insurance charges accurately, especially for non-smokers or low-risk individuals, but struggles with high charges due to outliers. With an R² of 0.78, the model explains most variation, and incorporating non-linear terms or robust methods could improve predictions for high-risk cases.

```
# Lasso Regression (Regularization)
lasso = Lasso(alpha=0.1)
lasso.fit(X_train, y_train)
y_lasso_pred = lasso.predict(X_test)
print("Lasso Regression R-squared:", r2_score(y_test, y_lasso_pred))
```

Lasso Regression R-squared: 0.7835913633542768

```
# Ridge Regression (Regularization)
ridge = Ridge(alpha=1.0)
ridge.fit(X_train, y_train)
y_ridge_pred = ridge.predict(X_test)
print("Ridge Regression R-squared:", r2_score(y_test, y_ridge_pred))
```

Ridge Regression R-squared: 0.7834446266673823

```
# Elastic Net Regression
from sklearn.linear_model import ElasticNet
elastic = ElasticNet(alpha=0.1, l1_ratio=0.5)  # l1_ratio=0.5 balances L1 and L2
elastic.fit(X_train, y_train)
```



```
# Elastic Net Regression
from sklearn.linear_model import ElasticNet
elastic = ElasticNet(alpha=0.1, l1_ratio=0.5)  # l1_ratio=0.5 balances L1 and L2
elastic.fit(X_train, y_train)
y_elastic_pred = elastic.predict(X_test)
print("Elastic Net R-squared:", r2_score(y_test, y_elastic_pred))
```

Elastic Net R-squared: 0.766737174732858

The performance metrics clearly show the impact of including additional predictors in the model. The Simple Linear Regression (SLR) model using only one predictor achieves a very low R-squared of 0.12 and a high mean squared error of 135,983,957.48, indicating poor predictive power. In contrast, the Multiple Linear Regression (MLR) model, which incorporates multiple predictors, achieves a much higher R-squared of 0.78 and a substantially lower mean squared error of 33,596,915.85. This demonstrates that adding relevant features significantly improves the model's ability to explain the variation in insurance charges and reduces prediction errors.