**Name :** Kaushik Kotian     **Roll no.** 30     **Div :** D15B     **Batch:**B

# DMBI Lab

# EXPERIMENT NO. 7

**AIM : To implement a regression model using Rapid Miner and Python.**
 1. Preprocess data. Split data into train and test set
 2. Build Regression model using inbuilt library function on training data
 3. Calculate metrics based on test data using inbuilt function
 4. Build a Regression model using a function defined by the student.
 5. Calculate metrics based on test data using inbuilt function
 6. Compare the results of all three ways of implementation.(Rapid Miner, Python Library)

**Theory :**
To implement a regression model using RapidMiner and Python, you can leverage both user-defined functions and built-in functions. Here's a general outline of how you can approach this:

**Data Preparation:**

- Load your dataset into RapidMiner for preprocessing. This may involve cleaning missing values, handling categorical variables, and scaling numeric features.
- Export the preprocessed data from RapidMiner to a format compatible with Python, such as CSV or Excel.

**Regression Model Building in RapidMiner:**

- Use RapidMiner's built-in operators for regression analysis, such as Linear Regression, Decision Tree Regression, or Support Vector Regression, depending on your data and problem.
- Configure the parameters of the regression model within RapidMiner, such as selecting input variables, setting regularization options, and specifying the target variable.

**Exporting the Model:**

- Once you have trained and validated your regression model in RapidMiner, export the model as a file (e.g., PMML format) that can be loaded into Python.
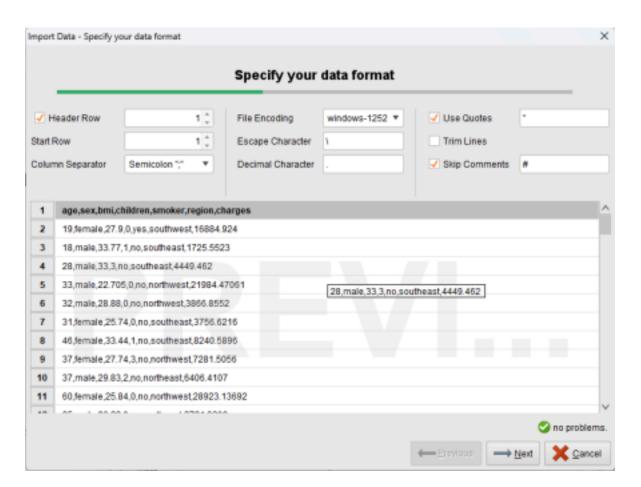
**Loading the Model in Python:**

- Use Python libraries such as pandas to load the preprocessed data and sklearn to load the exported regression model from RapidMiner.
- If needed, define custom functions in Python for any specific data transformations or model evaluation metrics that are not directly available in RapidMiner.

**Prediction and Evaluation:**

- Use the loaded regression model in Python to make predictions on new data or evaluate its performance on a test dataset.
- Implement evaluation metrics such as Mean Squared Error (MSE), R-squared, or others to assess the model's accuracy and reliability.

**IMPLEMENTATION USING RAPID MINER**

## Import Data - Format your columns.
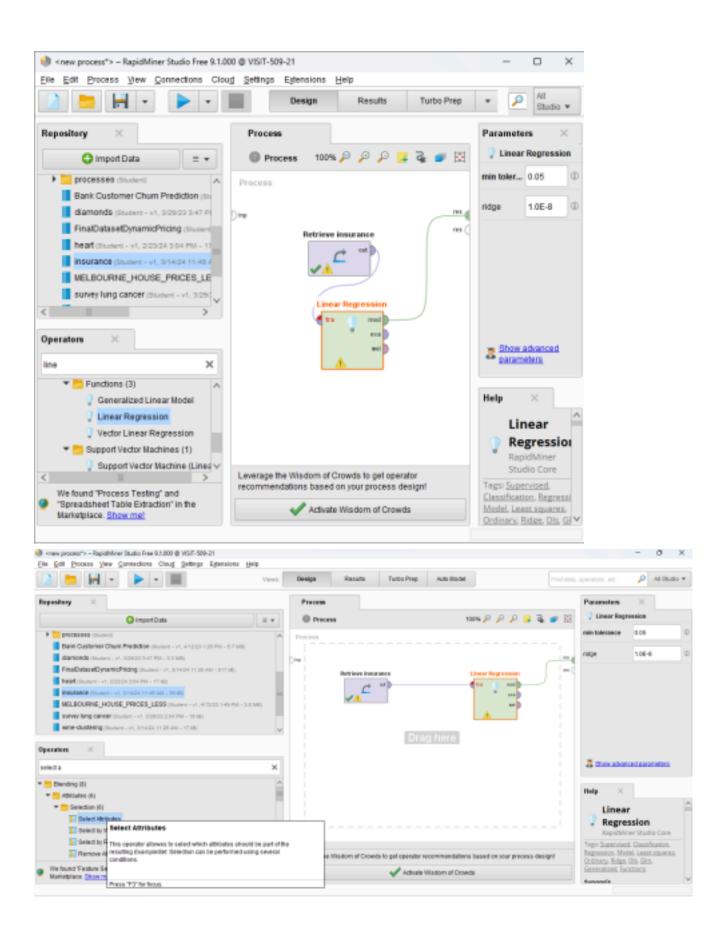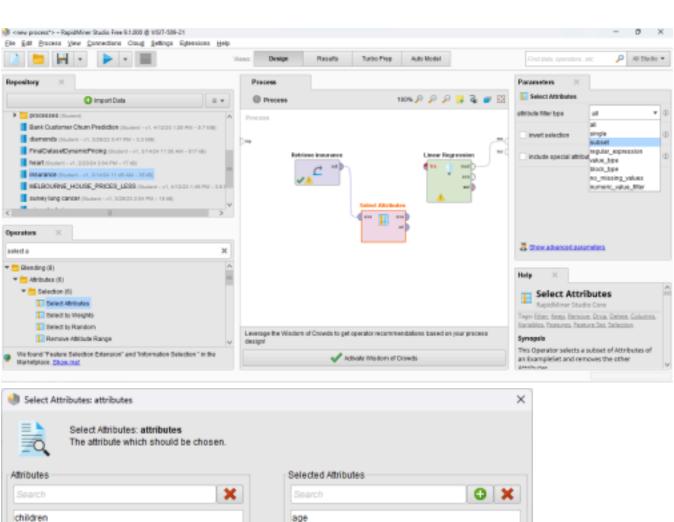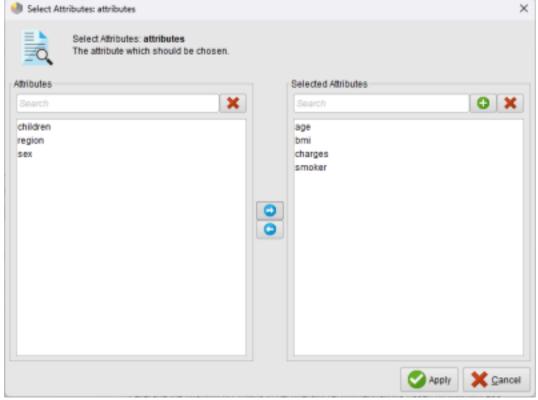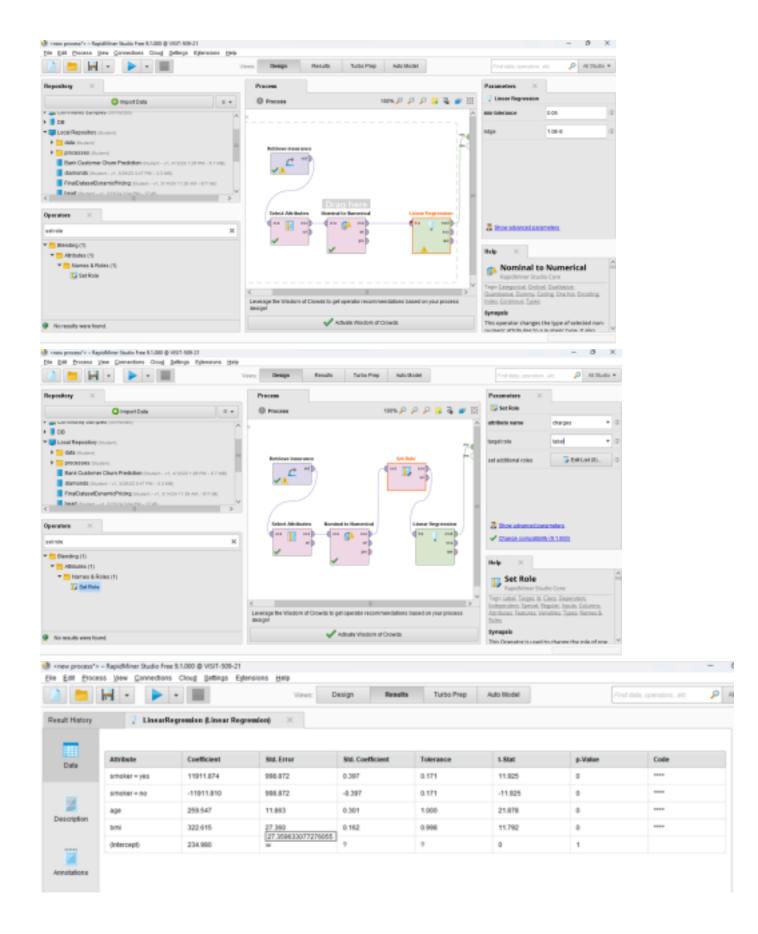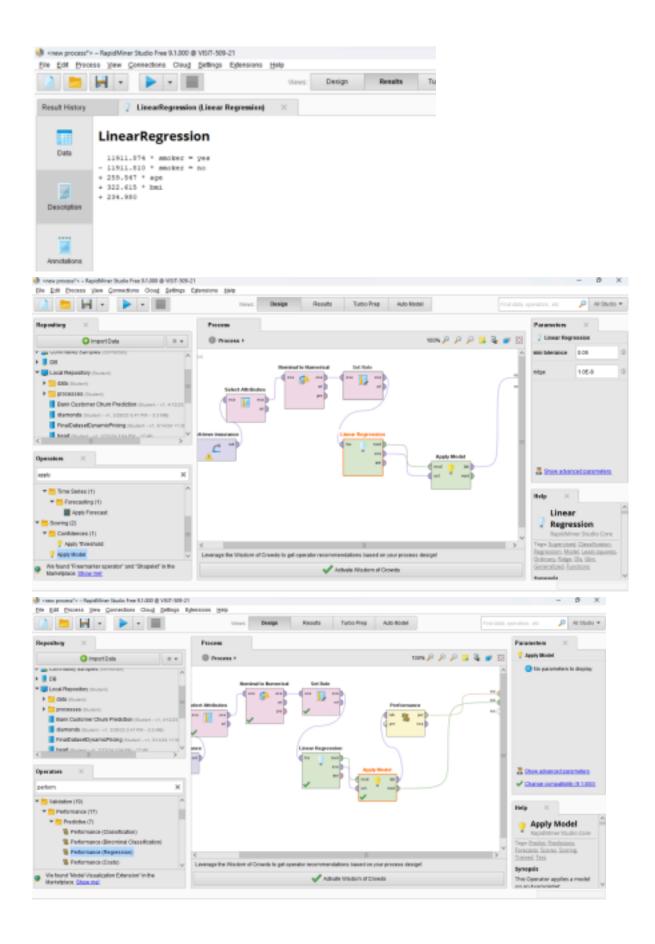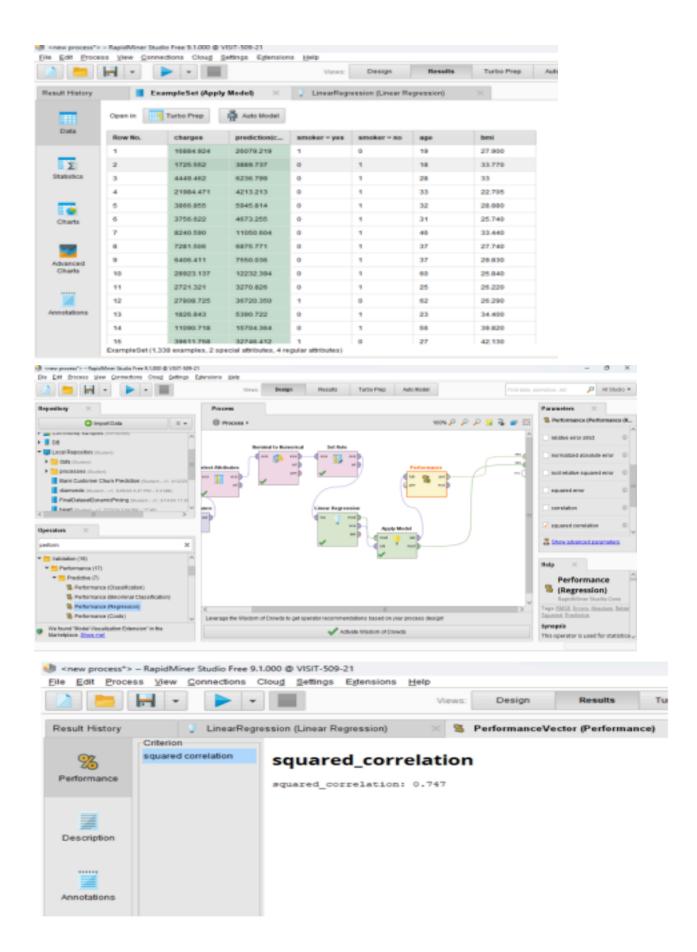
# Format your columns.

Date format: `Enter value...` ▼          ☐ Replace errors with missing values ⓘ

| | age<br>integer | sex<br>polynominal | bmi<br>real | children<br>integer | smoker<br>polynominal | region<br>polynominal | c<br>n |
|---|---|---|---|---|---|---|---|
| 1 | 19 | female | 27.900 | 0 | yes | southwest | |
| 2 | 18 | male | 33.770 | 1 | no | southeast | |
| 3 | 28 | male | 33.000 | 3 | no | southeast | |
| 4 | 33 | male | 22.705 | 0 | no | northwest | |
| 5 | 32 | male | 28.880 | 0 | no | northwest | |
| 6 | 31 | female | 25.740 | 0 | no | southeast | |
| 7 | 46 | female | 33.440 | 1 | no | southeast | |
| 8 | 37 | female | 27.740 | 3 | no | northwest | |
| 9 | 37 | male | 29.830 | 2 | no | northeast | |
| 10 | 60 | female | 25.840 | 0 | no | northwest | |
| 11 | 25 | male | 26.220 | 0 | no | northeast | |

✔ no problems.

⟵ Previous        ⟶ Next        ✖ Cancel

---

## \<new process*\> – RapidMiner Studio Free 9.1.000 @ VISIT-509-21

File  Edit  Process  View  Connections  Cloud  Settings  Extensions  Help

**Design**    Results    Turbo Prep    ▼    🔍    All Studio ▼

### Repository

🟢 Import Data        ≡ ▼

▶ 📁 processes (Student)
📘 Bank Customer Churn
📘 diamonds (Student - v1,
📘 FinalDatasetDynamicP
📘 heart (Student - v1, 2/232
📘 insurance (Student - v1,
📘 MELBOURNE_HOUSE
📘 survey lung cancer (Stu

### Operators

line                          ✖

▼ 📁 Functions (3)
  💡 Generalized Line
  💡 Linear Regres
  💡 Vector Linear R
▼ 📁 Support Vector Ma

We found "Process Testing" and "Spreadsheet Table Extraction" in the Marketplace. Show me!

### Process

《 100% 🔍 🔍 🔍 📄 📑 🗂 🔲

Process

Retrieve insurance

[Retrieve insurance operator]  out

Drag here

Activate Wisdom of Crowds

### Parameters

↪ Retrieve insurance (Retr...

repository e... /Loc  📁  ⓘ

Show advanced parameters

### Help

**Linear Regression**
Linear regression.

Press "F3" for focus.

...recommendations based on your process design!

Tags: Supervised, Classification, Regression, Model, Least squares, Ordinary, Ridge, Ols, Glm,

## LinearRegression (Linear Regression)

| Attribute | Coefficient | Std. Error | Std. Coefficient | Tolerance | t-Stat | p-Value | Code |
|-----------|-------------|-----------|------------------|-----------|--------|---------|------|
| smoker = yes | 11911.874 | 998.872 | 0.397 | 0.171 | 11.925 | 0 | **** |
| smoker = no | -11911.810 | 998.872 | -0.397 | 0.171 | -11.925 | 0 | **** |
| age | 259.547 | 11.863 | 0.301 | 1.000 | 21.878 | 0 | **** |
| bmi | 322.615 | 27.380 / 27.359633077276055 | 0.162 | 0.998 | 11.792 | 0 | **** |
| (Intercept) | 234.980 | | ? | ? | 0 | 1 | |

LinearRegression

```
11911.874 * smoker = yes
- 11911.810 * smoker = no
+ 259.547 * age
+ 322.615 * bmi
+ 234.980
```

RapidMiner Studio screenshots:

**Screenshot 1 — ExampleSet (Apply Model) Results**

<new process*> – RapidMiner Studio Free 9.1.000 @ VISIT-509-21

File  Edit  Process  View  Connections  Cloud  Settings  Extensions  Help

Views:  Design  Results  Turbo Prep  Auto

Result History  |  ExampleSet (Apply Model)  |  LinearRegression (Linear Regression)

Open in  Turbo Prep  Auto Model

| Row No. | charges | prediction(c... | smoker = yes | smoker = no | age | bmi |
|---|---|---|---|---|---|---|
| 1 | 16884.924 | 26079.219 | 1 | 0 | 19 | 27.900 |
| 2 | 1725.552 | 3889.737 | 0 | 1 | 18 | 33.770 |
| 3 | 4449.462 | 6236.799 | 0 | 1 | 28 | 33 |
| 4 | 21984.471 | 4213.213 | 0 | 1 | 33 | 22.705 |
| 5 | 3866.855 | 5945.814 | 0 | 1 | 32 | 28.880 |
| 6 | 3756.622 | 4673.255 | 0 | 1 | 31 | 25.740 |
| 7 | 8240.590 | 11050.604 | 0 | 1 | 46 | 33.440 |
| 8 | 7281.506 | 8875.771 | 0 | 1 | 37 | 27.740 |
| 9 | 6406.411 | 7560.036 | 0 | 1 | 37 | 29.830 |
| 10 | 28923.137 | 12232.394 | 0 | 1 | 60 | 25.840 |
| 11 | 2721.321 | 3270.826 | 0 | 1 | 25 | 26.220 |
| 12 | 27808.725 | 36720.350 | 1 | 0 | 62 | 26.290 |
| 13 | 1826.843 | 5390.722 | 0 | 1 | 23 | 34.400 |
| 14 | 11090.718 | 15704.364 | 0 | 1 | 56 | 39.820 |
| 15 | 39611.758 | 32748.412 | 1 | 0 | 27 | 42.130 |

ExampleSet (1,338 examples, 2 special attributes, 4 regular attributes)

**Screenshot 2 — Design view**

Parameters: Performance (Performance (R...

- relative error strict
- normalized absolute error
- root relative squared error
- squared error
- correlation
- ☑ squared correlation

Show advanced parameters

Help — Performance (Regression) — RapidMiner Studio Core

Operators: perform
- Validation (10)
- Performance (17)
  - Predictive (7)
    - Performance (Classification)
    - Performance (Binominal Classification)
    - Performance (Regression)
    - Performance (Costs)

Leverage the Wisdom of Crowds to get operator recommendations based on your process design.

Activate Wisdom of Crowds

**Screenshot 3 — PerformanceVector (Performance)**

Result History  |  LinearRegression (Linear Regression)  |  PerformanceVector (Performance)

Criterion: squared correlation

**squared_correlation**

squared_correlation: 0.747

## 2. IMPLEMENTATION USING BUILT IN FUNCTION :

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
data = pd.read_csv('wine-clustering.csv')
print(data.head())
```

```
   Alcohol  Malic_Acid  Ash  Ash_Alcanity  Magnesium  Total_Phenols  \
0    14.23        1.71  2.43          15.6        127           2.80
1    13.20        1.78  2.14          11.2        100           2.65
2    13.16        2.36  2.67          18.6        101           2.80
3    14.37        1.95  2.50          16.8        113           3.85
4    13.24        2.59  2.87          21.0        118           2.80

   Flavanoids  Nonflavanoid_Phenols  Proanthocyanins  Color_Intensity   Hue  \
0        3.06                  0.28             2.29             5.64  1.04
1        2.76                  0.26             1.28             4.38  1.05
2        3.24                  0.30             2.81             5.68  1.03
3        3.49                  0.24             2.18             7.80  0.86
4        2.69                  0.39             1.82             4.32  1.04

   OD280  Proline
0   3.92     1065
1   3.40     1050
2   3.17     1185
3   3.45     1480
4   2.93      735
```

```python
# Assuming 'X' contains your input features and 'y' contains the target variable
X = data[['Malic_Acid', 'Ash_Alcanity', 'Ash','Flavanoids','Color_Intensity']]
y = data['Alcohol']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = LinearRegression()
model.fit(X_train, y_train)
# Make predictions on the test data
y_pred = model.predict(X_test)
# Calculate evaluation metrics
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"Mean Squared Error (MSE): {mse}")
print(f"R-squared (R2): {r2}")
```

## 3. IMPLEMENTATION USING USER DEFINED FUNCTION :

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
```

```python
def train_regression_model(data, features, target, test_size=0.2, random_state=42):
    # Split the data into input features (X) and target variable (y)
    X = data[features]
    y = data[target]
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size, random_state=random_state)
    model = LinearRegression()
    model.fit(X_train, y_train)
    return model, X_test, y_test
```

```python
def evaluate_regression_model(model, X_test, y_test):
    # Make predictions on the test data
    y_pred = model.predict(X_test)
    # Calculate evaluation metrics
    mse = mean_squared_error(y_test, y_pred)
    r2 = r2_score(y_test, y_pred)
    return mse, r2
```

```python
data = pd.read_csv('wine-clustering.csv')
features = ['Malic_Acid', 'Ash_Alcanity', 'Ash','Flavanoids','Color_Intensity']
target = 'Alcohol'

model, X_test, y_test = train_regression_model(data, features, target)
mse, r2 = evaluate_regression_model(model, X_test, y_test)
print(f"Mean Squared Error (MSE): {mse}")
print(f"R-squared (R2): {r2}")
```

```
Mean Squared Error (MSE): 0.2650597239422124
R-squared (R2): 0.5560404872772855
```

**COMPARISON AND CONCLUSION :** Comparing the Mean Squared Error (MSE) across Python's user-defined functions, Python's built-in functions, and RapidMiner reveals varying levels of flexibility, complexity, and performance. Python with user-defined functions allows for fine-tuning and optimization, potentially leading to lower MSE. Python's built-in functions offer a balance between simplicity and performance. RapidMiner's MSE depends on the efficiency of its built-in operators and workflow design. The Python implementation offers flexibility and control for customized data processing and model training using libraries like pandas and scikit-learn.