**Sankeerth Sai Shabad**

**Introduction to Data Science Project-3**

**03/01/2022**

**PURPOSE OF THE PROBLEM:**

- To perform data cleaning and find and insert data by using data cleaning techniques and necessary steps to ensure the dataset is cleansed and performing SQL queries using MySQL Workbench.
- Importing cleansed MySQL files into data frames in Jupyter notebook in Google colab. And to perform and plot the linear regression model and independent models that have a high correlation with correlation analysis
- To provide the equation of each model and compute their R2 and MSE.

**Methodology:**

- COLLECTION OF DATA: All the data is collected from the dataset with their values from
  https://github.com/SankeerthShabad/IDS/blob/main/HW3/Life_Expectancy.csv
- OPERATIONS: Addressing missing values by using mean and performing required queries by using SQL in MySQL Workbench. And Importing cleansed data into Jupyter notebook in Google colab. And to perform and plot linear regression, correlation analysis by using Python.
- OBSERVATIONS: To find the values by comparing and effects caused by life expectancy and mortality rates. Negative correlation for life expectancy by social and economic factors. And the impact of schooling on the lifespan of humans and to observe the performance of the models.

**Results:**

- SQL: After performing SQL queries for data cleaning and operations required, using MySQL Workbench the following outputs are observed.
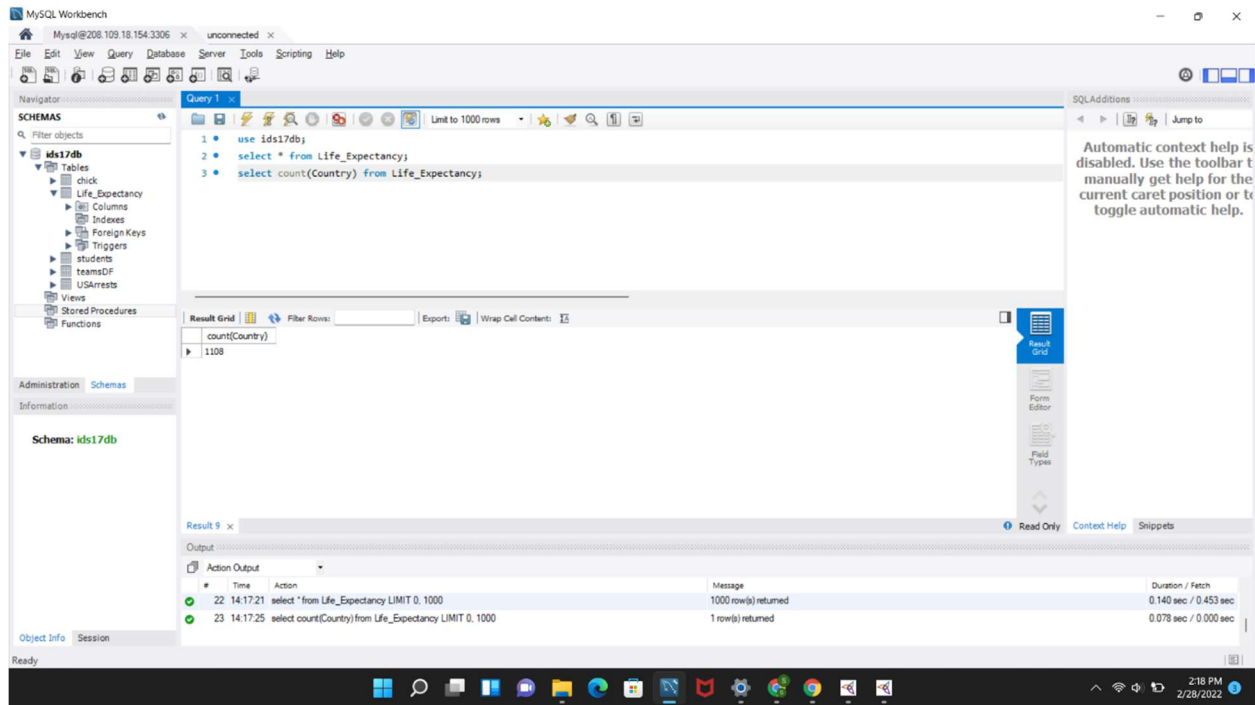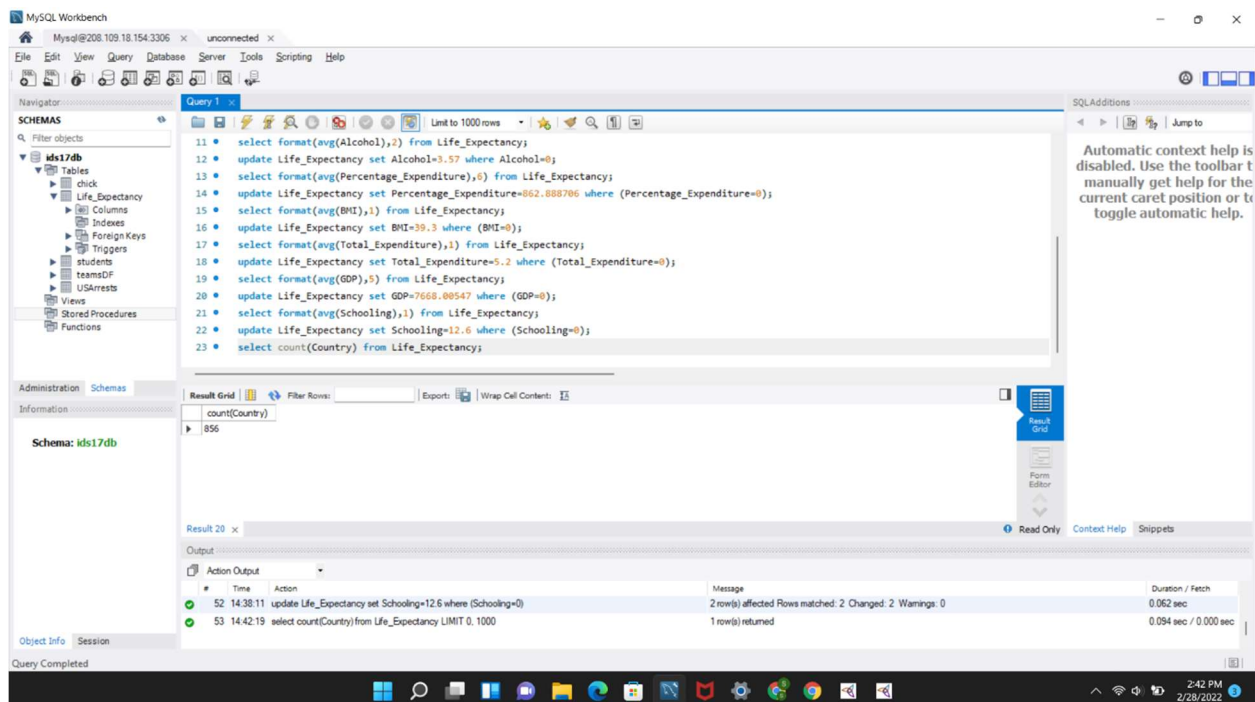
Fig: Total number of country counts



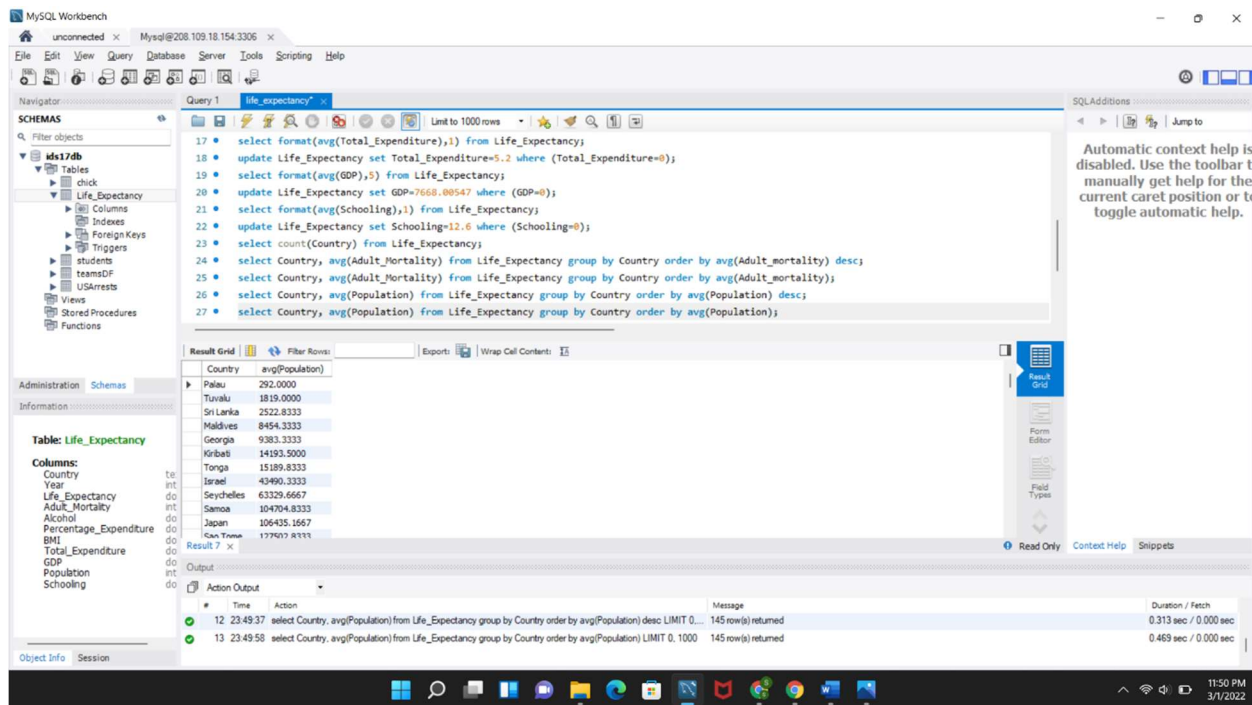Fig: Total count of countries after performing data cleaning techniques

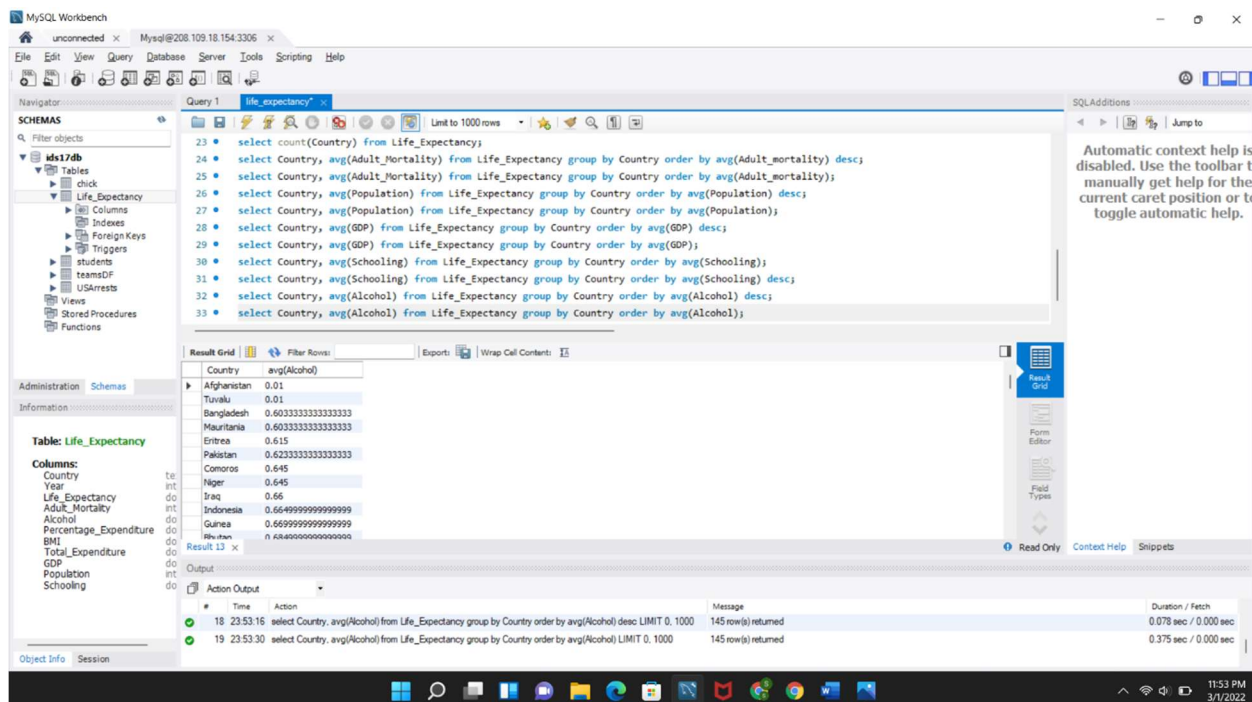Fig: List of countries with highest average mortality rates



Fig: List of countries with Lowest average mortality rates
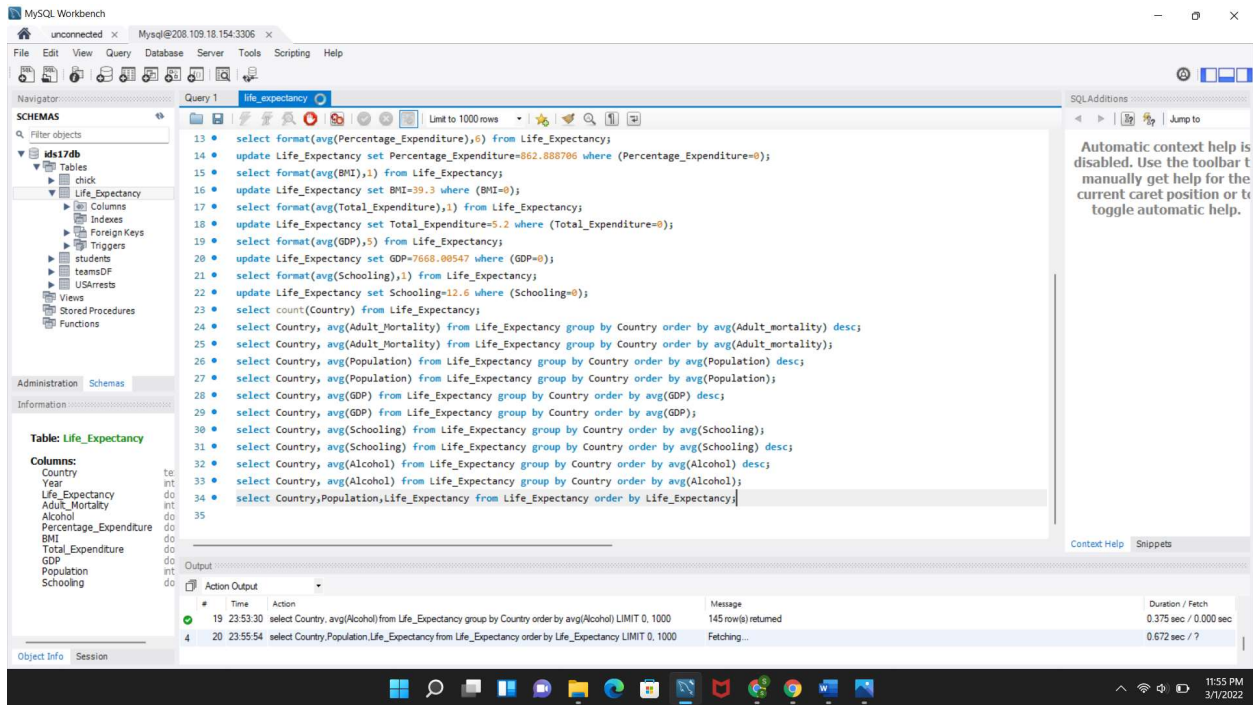
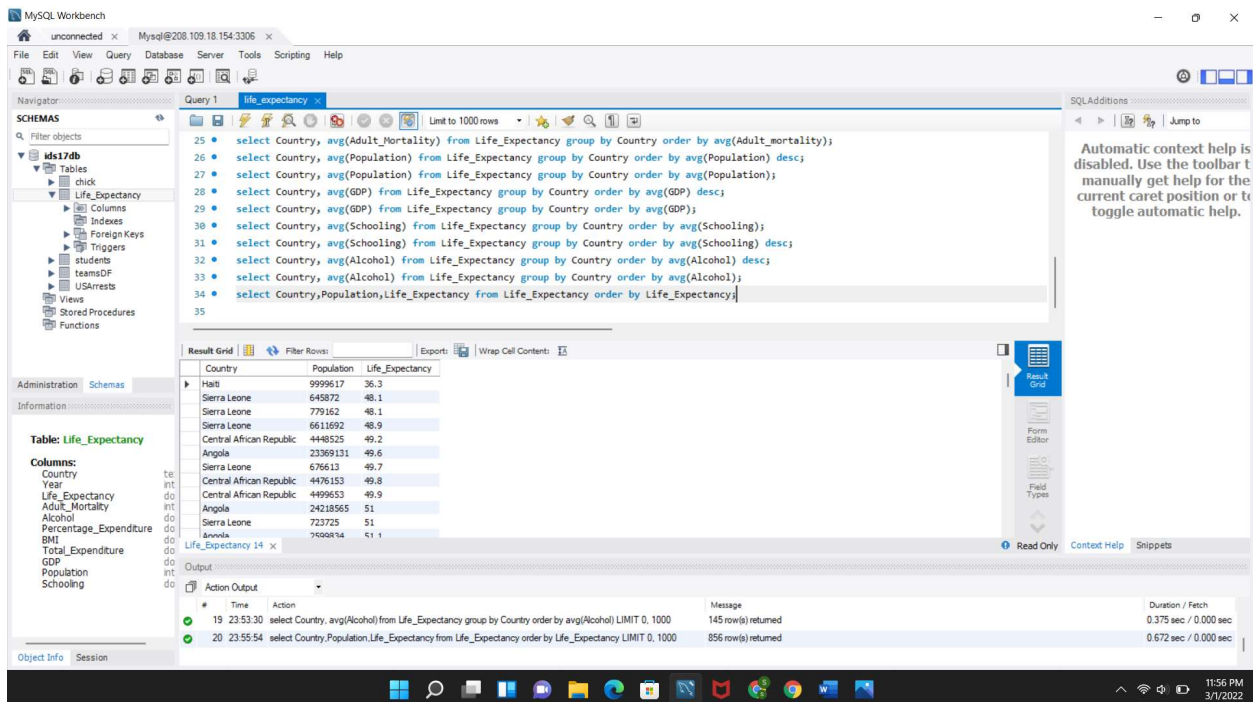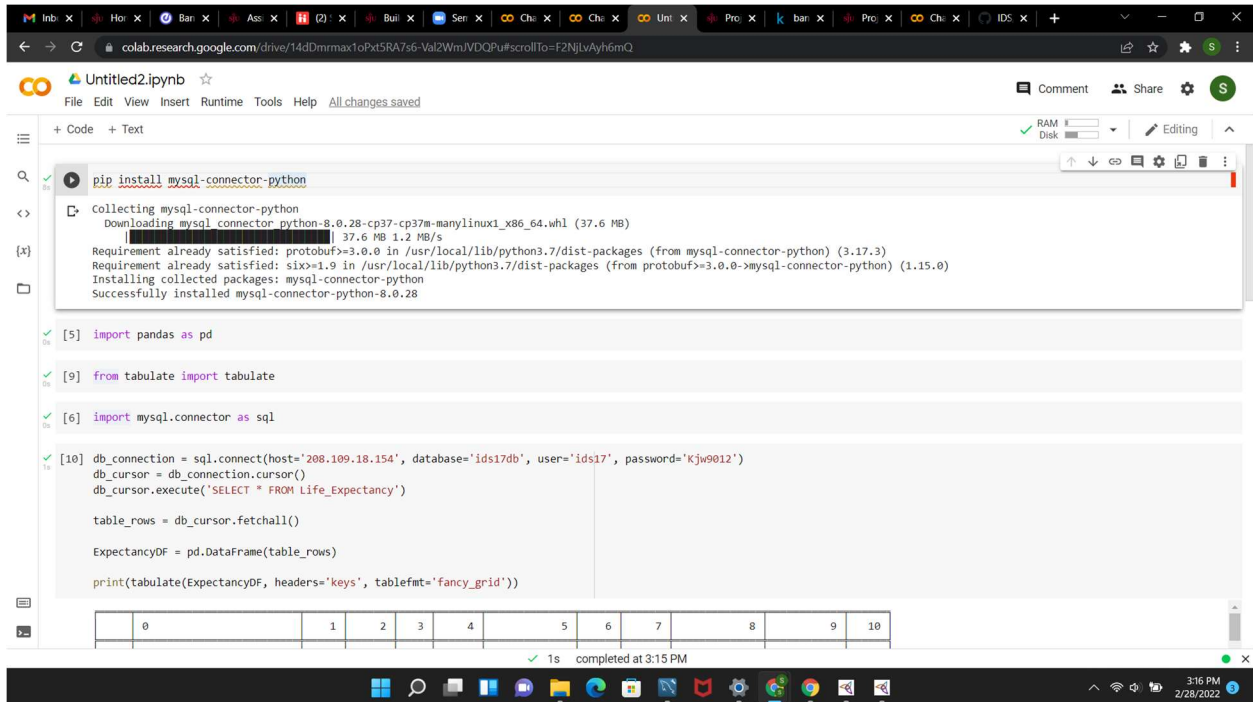Fig: Queries for the highest and lowest average for Countries and years



Fig: Comparing life expectancy for densely populated countries

**PYTHON:** After Importing cleansed MySQL files into data frames in Jupyter notebook in Google colab the plots, linear regression model, and independent models that have a high correlation with correlation analysis and computed R2 and MSE.



Fig: Importing cleansed data in Jupyter

Fig: calculating count, mean, std, min by using describe().



Fig: variance of expectancy

```python
from sklearn.metrics import mean_squared_error
y_true=[3,-0.5,2,7]
y_pred=[2.5,0.0,2,8]
mean_squared_error(y_true, y_pred)
```

```
0.375
```

```python
X = ExpectancyDF.loc[:,'Life_Expectancy'].values.reshape(-1,1)

Y = ExpectancyDF.loc[:,'Schooling'].values.reshape(-1,1)

LRmodel = LinearRegression()

LRmodel.fit(X, Y)

Y_pred = LRmodel.predict(X)

plt.scatter(X, Y)
plt.plot(X, Y_pred, color='red')

plt.show()
```

Fig: plotting linear regression model



```python
plt.show()
```

```python
print("The slope: ", LRmodel.coef_)
print("The intercept: ", LRmodel.intercept_)
```

```
The slope:  [[0.26892275]]
The intercept:  [-6.32507767]
```

```python
from sklearn.metrics import mean_squared_error, r2_score

print("MSE: ", mean_squared_error(Y, Y_pred))
print("R2: ", r2_score(Y, Y_pred))
```

```
MSE:  3.128838215956491
R2:  0.645797347147856
```
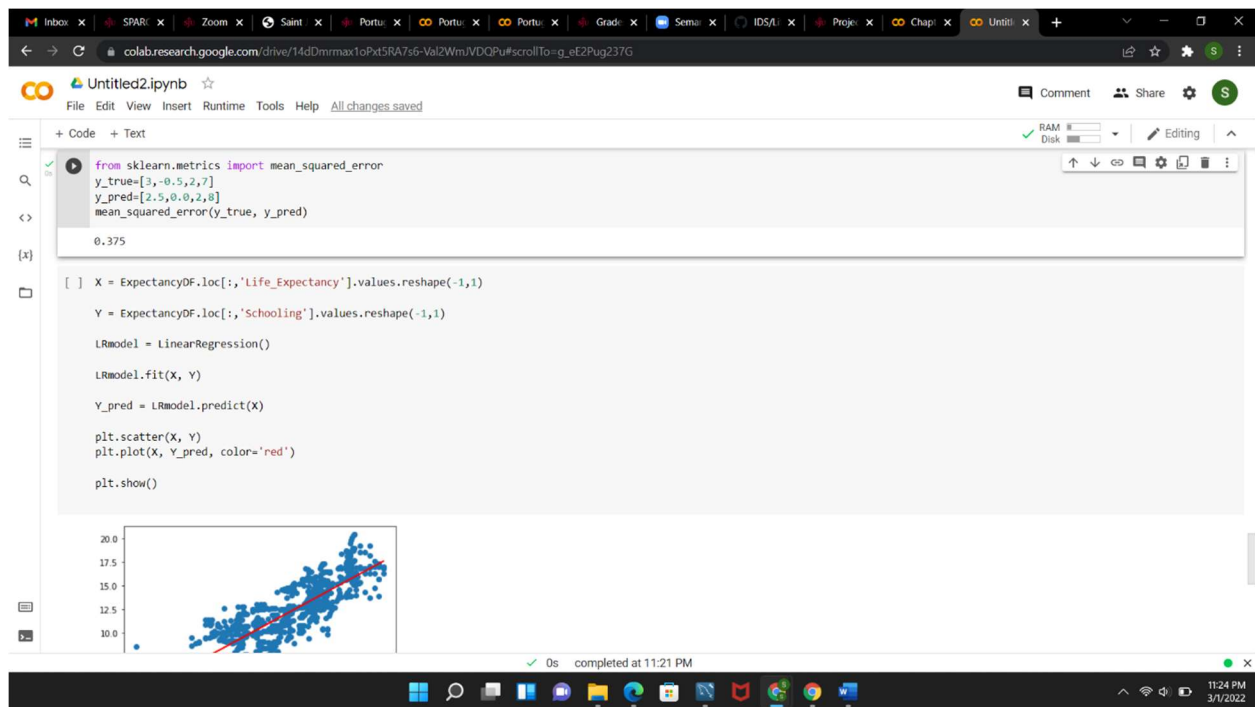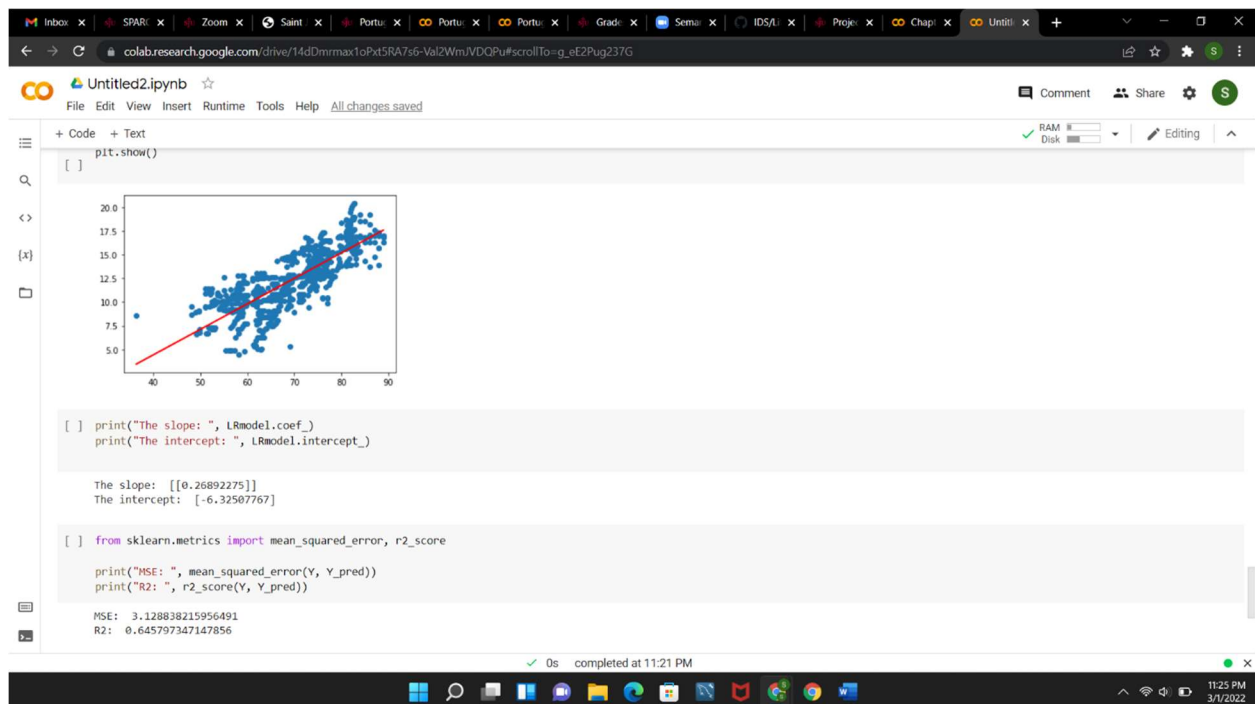
Fig: Printing slope, intercept, MSE, R2 score

**Conclusion:** Several linear regression models for life expectancy as found in correlation analysis and I computed R2 and MSE values. After summarizing Regression model performs best.