| **Prelim Examination** | |
|---|---|
| Classification and Prediction Analysis Report | |
| **Course Code:** CPE019 | **Program:** BSCpE |
| **Course Title:** Emerging Technologies 2 | **Date Performed:** Feb 28, 2024 |
| **Section:** CpE32S3 | **Date Submitted:** Mar 6, 2024 |
| **Name:**<br>Maglalang, Charles Lester<br>Nicolas, Sean Julian | **Instructor:** Engr. Roman Richard |

**Import Libraries and Creating Functions**

In this stage we imported all of the libraries in a single cell block so that we can import all of the needed libraries in one execution. Creating functions including the lines of code that we often use makes it convenient and fast. Two functions included in the picture below are convertStrNum() that convert words like "yes" and "no" into numbers, the other function is importdata() which is used to easily import the dataset. In the code below, we also cloned the github path in order to easily access the dataset file.

```
[1] import numpy as np
    from sklearn.linear_model import LinearRegression, LogisticRegression
    from sklearn.preprocessing import PolynomialFeatures, StandardScaler, LabelEncoder
    from sklearn.model_selection import train_test_split
    from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
    from sklearn.metrics import r2_score as score
    from sklearn.metrics import accuracy_score, mean_squared_error, classification_report
    from sklearn.tree import plot_tree, DecisionTreeClassifier
    import matplotlib.pyplot as plt
    from mpl_toolkits.mplot3d import Axes3D
    import pandas as pd
    import seaborn as sns

    !git clone "https://github.com/SeanNicolas/CPE019.git"

    Cloning into 'CPE019'...
    remote: Enumerating objects: 50, done.
    remote: Counting objects: 100% (50/50), done.
    remote: Compressing objects: 100% (44/44), done.
    remote: Total 50 (delta 12), reused 0 (delta 0), pack-reused 0
    Receiving objects: 100% (50/50), 1.14 MiB | 8.64 MiB/s, done.
    Resolving deltas: 100% (12/12), done.

[2] def convertStrNum(col, dataset):
        le = LabelEncoder()
        dataset[col]= le.fit_transform(dataset[col])

[3] def importdata(path):
        dataset = pd.read_csv(path)
        return dataset

[4] ins = importdata('/content/CPE019/Prelim Exam/insurance.csv')

[5] convertStrNum("smoker", ins)
    convertStrNum("sex", ins)
    convertStrNum("region", ins)
```

*Importing Libraries and Creating Function*

## Identifying Correlation

Using heatmap from the seaborn library, we can check the correlation between the different variables inside the database. This is helpful to identify how each variable affects one another and what variable are they most correlated with.
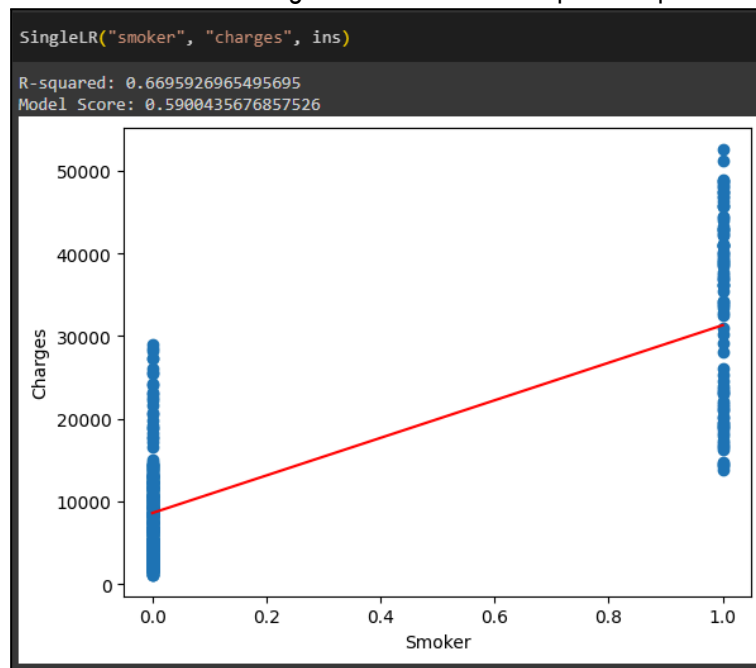
```
sns.heatmap(ins.corr())
```
```
<Axes: >
```



*Using Heatmap to Check Correlations*

## Singular Linear Regression

In Singular Linear Regression, the libraries that we used are Pandas and Numpy for data manipulation and reshaping arrays. Matplotlib for plotting and visualization, with sklearn.linear_model it is used for linear regression modeling. The dataset that we chose for this model is Insurance.csv, for X we used 'Smoker' and 'Charges' for our Y since they have a high relationship and fitted in this model.

It is shown here the visualization of the regression line in scatter plot and predicted values.

```
SingleLR("smoker", "charges", ins)

R-squared: 0.6695926965495695
Model Score: 0.5900435676857526
```
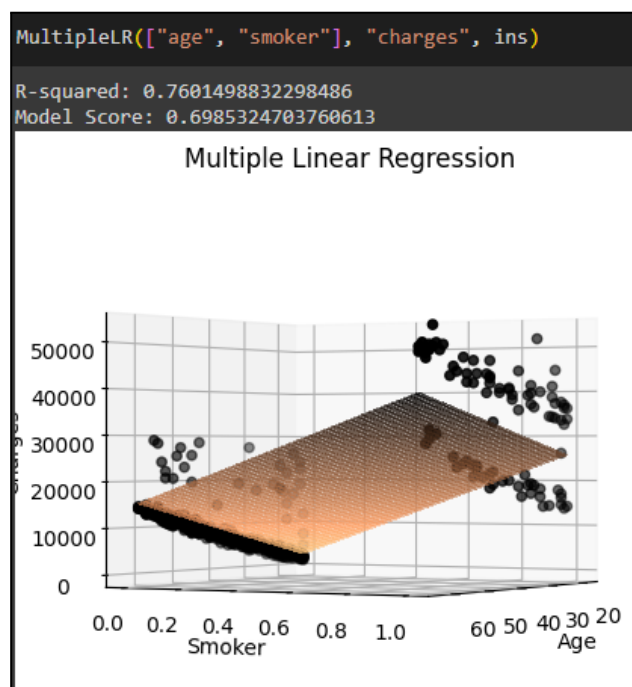


*Linear Regression*

Looking at the figure above, we are able to see a linear regression between the smoker and charges variables. The connection between smoker and charges offered the highest accuracy after trying different combinations of variables

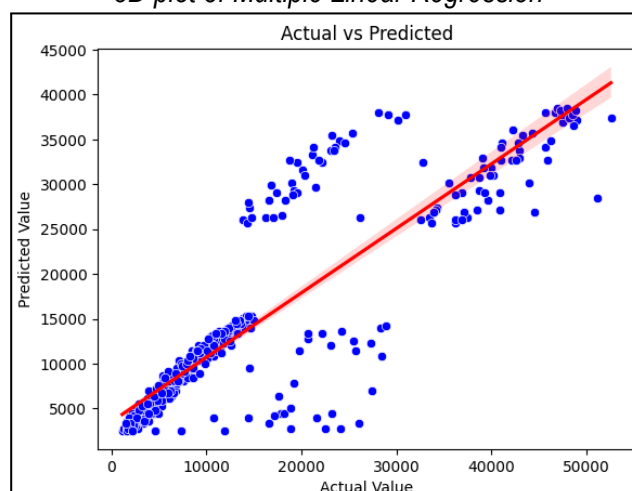**Multiple Linear Regression**

Just like the Singular Linear Regression, we used the same libraries but we added Seaborn, it is the same with matplotlib but it has a high-level interface. We split the dataset into training and testing sets, incorporating multiple features into a data frame with Age, Smoker, Charges columns.

It shows here the visualization of Multiple Linear Regression using seaborn for 3d representation. Also, it shows here the actual and predicted values from the variables that we used, model score is included to show the accuracy of the model.

Upon checking the different combinations of variables, using age and smoker to predict charges had the highest accuracy

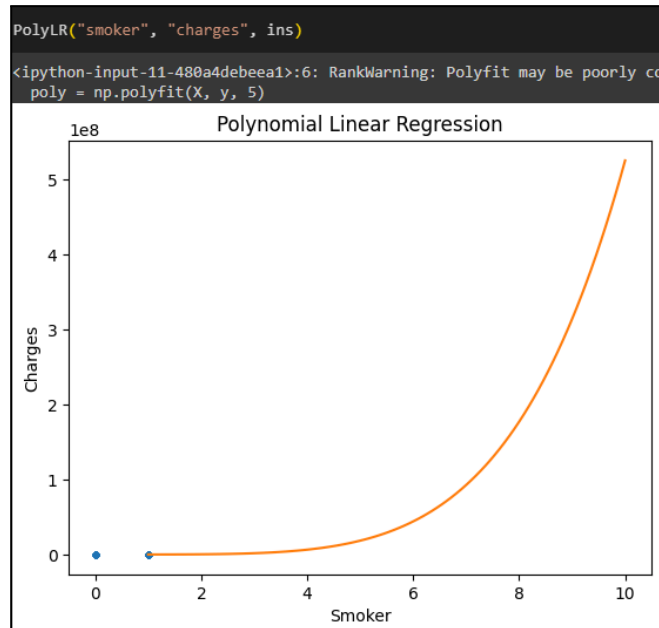

*3D plot of Multiple Linear Regression*



*Prediction vs Actual*

**Polynomial Linear Regression**

Compared to the other regression models, this regression needs to view the scatter plot to identify what polynomial feature will be used. In this model, we used the 'charges' and 'smoker' columns; these values represent the prediction of the polynomial curve for a given range of feature values.

It shows here the polynomial linear regression visualization, the dots represent the original data and the predicted values.
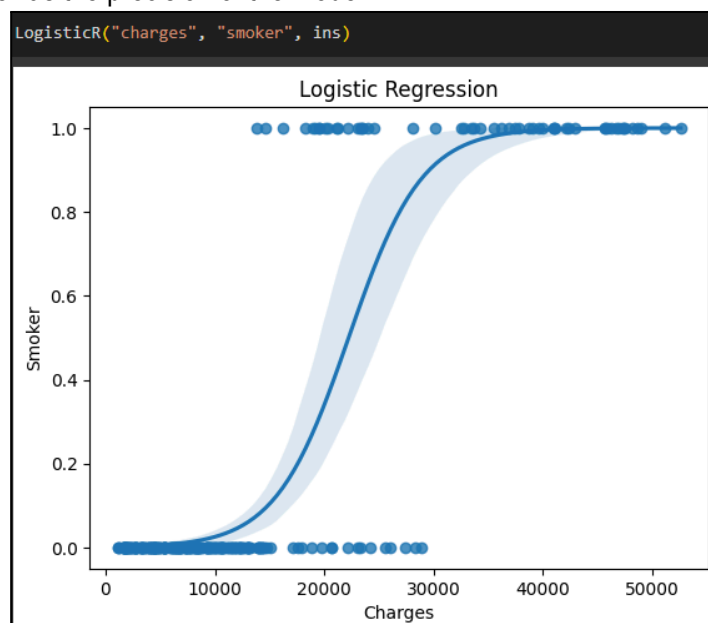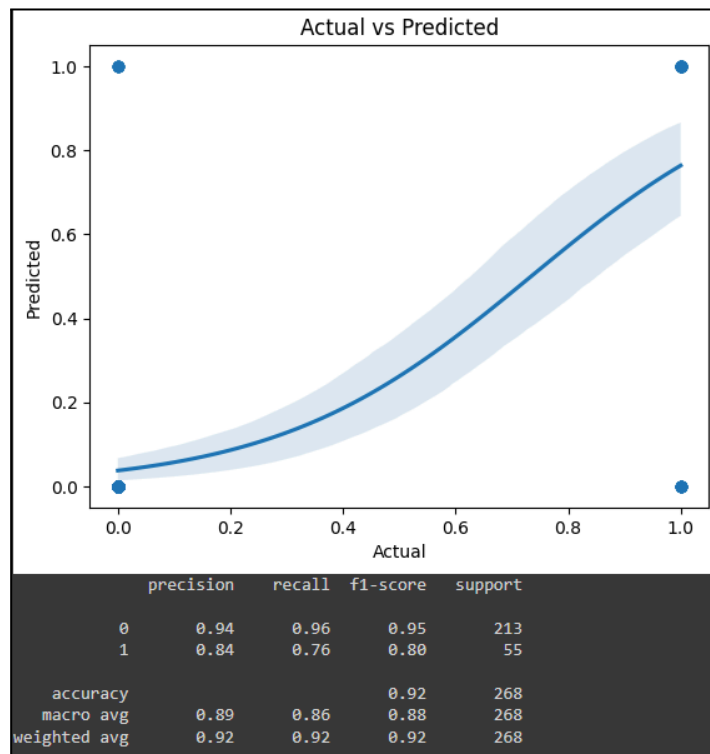


*Polynomial Linear Regression*

**Logistic Regression Model**

In a logistic regression model, libraries that need to be used are Pandas for data manipulation, Matplotlib for creating plots and visualization, and sklearn.linear_model.LogisticsRegression to build logistic regression models. Insurance.csv is the chosen data for this model since it fits in the regression line. 'charges' and 'smoker' columns will be the variables for our x and y.

It shows here the visualization of the logistic regression model, and the actual and predicted values are included, as well as the precision of the model.



Logistic Regression

*Predicted vs Actual and Evaluation of Model*

## Decision Tree

In the decision tree model, the same libraries are used, just like in the previous models. The data set was organized and separated into target values using train_test_split. A decision tree classifier was created using the entropy criterion to determine the best split at each node. The classifier trained on the training data and made predictions based on entropy.

It shows here the visualization of the decision tree and the accuracy of the model.



*Evaluation of Model*

charges <= 14997.505
entropy = 0.728
samples = 936
value = [746, 190]

age <= 18.5
entropy = 0.029
samples = 681
value = [679, 2]

charges <= 33611.133
entropy = 0.831
samples = 255
value = [67, 188]

## Preview of Decision Tree

## Full Decision Tree

## Random Forest

In creating a random forest, we found out that we can identify how many smokers there are in a region based on age, sex, BMI, and charge columns. The classifier used a default of 100 trees in the forest to predict test data. We created two variables, X and Y, where X represents the original data frame and Y represents the new data frame.

It shows here the visualization of random forests and the accuracy of the model. The predicted values for smokers in every region are also shown here.

```
RandomF(["age", "sex", "bmi", "charges"], "smoker", ins)

              precision    recall  f1-score   support

          0       0.94      0.98      0.96       318
          1       0.92      0.77      0.84        84

   accuracy                           0.94       402
  macro avg       0.93      0.88      0.90       402
weighted avg       0.94      0.94      0.94       402

Accuracy: 0.9378109452736318
```
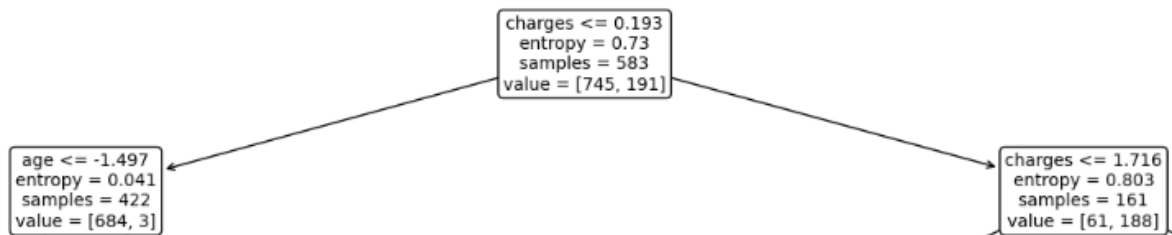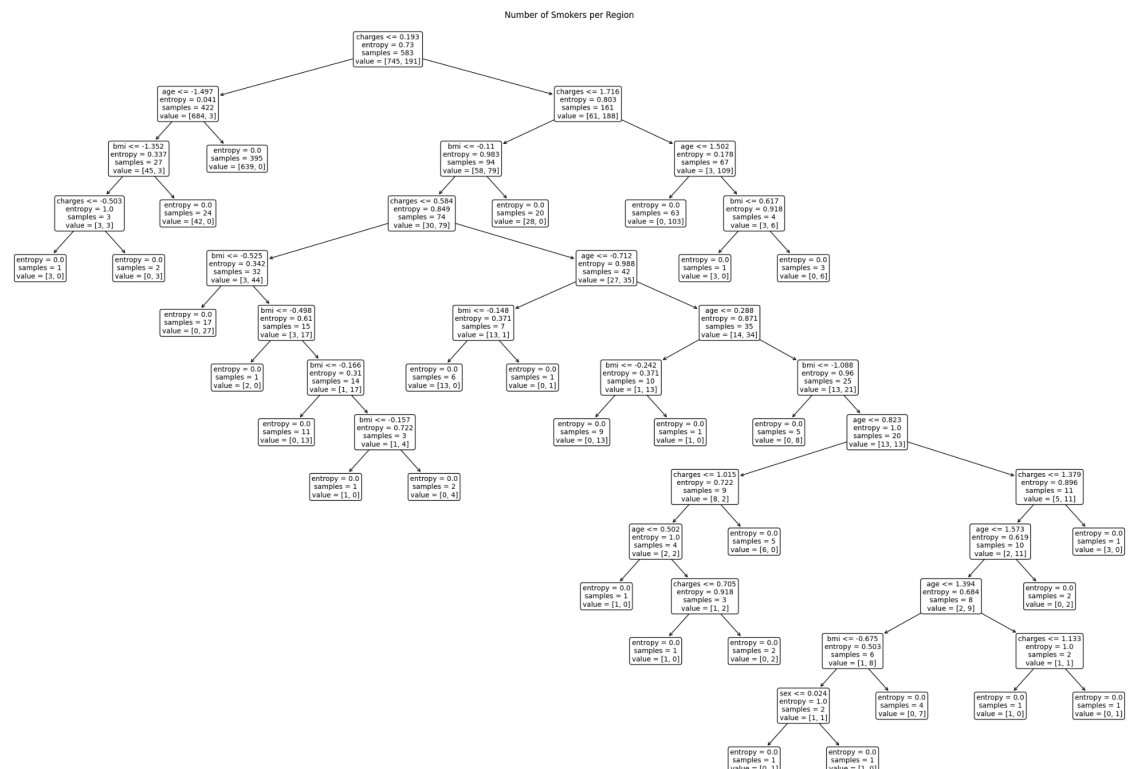


Preview



Number of Smokers per Region

Full View