

Medical Insurance Cost Prediction Project

Explore the future of healthcare with our cutting-edge machine learning project, which predicts medical insurance costs with precision and empowers informed decisions for individuals and providers.

Improting The libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn import metrics
```

Data Collection & Analysis

```
In [2]: # loading the data from csv file to a Pandas DataFrame
insurance_data = pd.read_csv('C:\Users\hvh\OneDrive\Desktop\insurance datafile.csv')

In [3]: insurance_data.head()

Out[4]:
   age  sex  bmi  children  smoker  region  charges
0   19  female  27.900    0    yes  southwest  16884.92400
1   18  male  33.770    1    no  southeast  1725.55230
2   28  male  33.000    3    no  southeast  4449.46200
3   33  male  22.705    0    no  northwest  21984.47061
4   32  male  28.880    0    no  northwest  3866.85520

...

1333  50  male  30.970    3    no  northwest  10600.54830
1334  18  female  31.920    0    no  northeast  2205.98080
1335  18  female  36.850    0    no  southeast  1629.83350
1336  21  female  25.800    0    no  southwest  2007.94500
1337  61  female  29.070    0    yes  northwest  29141.36030

1338 rows x 7 columns

In [5]: #original Insurance Dataset
insurance_data

Out[5]:
   age  sex  bmi  children  smoker  region  charges
0   19  female  27.900    0    yes  southwest  16884.92400
1   18  male  33.770    1    no  southeast  1725.55230
2   28  male  33.000    3    no  southeast  4449.46200
3   33  male  22.705    0    no  northwest  21984.47061
4   32  male  28.880    0    no  northwest  3866.85520

...

1333  50  male  30.970    3    no  northwest  10600.54830
1334  18  female  31.920    0    no  northeast  2205.98080
1335  18  female  36.850    0    no  southeast  1629.83350
1336  21  female  25.800    0    no  southwest  2007.94500
1337  61  female  29.070    0    yes  northwest  29141.36030

1338 rows x 7 columns

In [6]: # number of rows and columns
insurance_data.shape

Out[6]:
(1338, 7)

In [7]: # Statistical Analysis of our Dataset
insurance_data.std()

In [8]: C:\Users\hvh\AppData\Local\Temp\ipykernel_13789\93751182.py:1: FutureWarning: The default value of numeric_only in DataFrame.std is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.
insurance_data.std()

Out[8]:
age          14.049969
bmi          6.098187
children     1.286492
charges    12110.01237
dtype: float64

In [9]: # getting some informations about the dataset
insurance_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
 #   Column  Non-Null Count  Dtype
---  --
 0   age     1338 non-null    int64
 1   sex     1338 non-null    object
 2   bmi     1338 non-null    float64
 3   children 1338 non-null    int64
 4   smoker  1338 non-null    object
 5   region  1338 non-null    object
 6   charges 1338 non-null    float64
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB

In [10]: # Categorical Features:
#sex
#smoker
#region

In [11]: # checking for missing values
insurance_data.isnull().sum()

Out[11]:
age      0
sex      0
bmi      0
children 0
smoker   0
region   0
charges  0
dtype: int64
```

Data Analysis

```
In [12]: # statistical Measures of the dataset
insurance_data.describe()

Out[12]:
   age      bmi  children  charges
count  1338.000000  1338.000000  1338.000000  1338.000000
mean    29.207025    30.663397    1.094918    13270.422265
std     14.049960    6.098187    1.205493    11210.01237
min     18.000000    15.960000    0.000000    1121.873900
25%     27.000000    26.296250    0.000000    4740.287150
50%     30.000000    30.400000    1.000000    9382.033000
75%     51.000000    34.693750    2.000000    16639.912515
max     64.000000    53.130000    5.000000    63770.428010

In [13]: # distribution of age value
sns.set()
plt.figure(figsize=(6,6))
sns.distplot(insurance_data['age'])
plt.title('Age Distribution')
plt.show()

C:\Users\hvh\AppData\Local\Temp\ipykernel_11789\363492332.py:4: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/d644147ed2974457a06372759bbe5751
sns.distplot(insurance_data['age'])

Age Distribution
Density
0.040
0.035
0.030
0.025
0.020
0.015
0.010
0.005
0.000
10 20 30 40 50 60 70
age
```

```
In [14]: # Gender column
plt.figure(figsize=(6,6))
sns.countplot(x='sex', data=insurance_data)
plt.title('Sex Distribution')
plt.show()

Sex Distribution
count
700
600
500
400
300
200
100
0
female male
sex
```

```
In [15]: insurance_data['sex'].value_counts()

male      676
female    662
Name: sex, dtype: int64

In [16]: # bmi distribution
plt.figure(figsize=(6,6))
sns.distplot(insurance_data['bmi'])
plt.title('BMI Distribution')
plt.show()

C:\Users\hvh\AppData\Local\Temp\ipykernel_11789\1926795408.py:3: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/d644147ed2974457a06372759bbe5751
sns.distplot(insurance_data['bmi'])

BMI Distribution
Density
0.07
0.06
0.05
0.04
0.03
0.02
0.01
0.00
10 20 30 40 50
bmi
```

Normal BMI Range --> 18.5 to 24.9

```
In [17]: # children column
plt.figure(figsize=(6,6))
sns.countplot(x='children', data=insurance_data)
plt.title('Children')
plt.show()

Children
count
600
500
400
300
200
100
0
0 1 2 3 4 5
children
```

```
In [18]: insurance_data['children'].value_counts()

0      574
1      324
2      240
3      157
4       25
5       16
Name: children, dtype: int64

In [19]: # smoker column
plt.figure(figsize=(6,6))
sns.countplot(x='smoker', data=insurance_data)
plt.title('smoker')
plt.show()

smoker
count
1000
800
600
400
200
0
yes no
smoker
```

```
In [20]: # region column
plt.figure(figsize=(6,6))
sns.countplot(x='region', data=insurance_data)
plt.title('region')
plt.show()

region
count
350
300
250
200
150
100
50
0
southwest southeast northwest northeast
region
```

```
In [21]: insurance_data['region'].value_counts()

southeast  364
southwest  325
northwest  325
northeast  324
Name: region, dtype: int64

In [22]: # distribution of charges value
plt.figure(figsize=(6,6))
sns.distplot(insurance_data['charges'])
plt.title('Charges Distribution')
plt.show()

C:\Users\hvh\AppData\Local\Temp\ipykernel_11789\3971177822.py:3: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/d644147ed2974457a06372759bbe5751
sns.distplot(insurance_data['charges'])

Charges Distribution
Density
7
6
5
4
3
2
1
0
-10000 0 10000 20000 30000 40000 50000 60000 70000
charges
```

Data Pre-Processing

```
Encoding the categorical features

In [23]: # encoding sex column
insurance_data.replace({'sex':{'male':0,'female':1}}, inplace=True)

3 # encoding 'smoker' column
insurance_data.replace({'smoker':{'yes':0,'no':1}}, inplace=True)

# encoding 'region' column
insurance_data.replace({'region':{'southeast':0,'southwest':1,'northeast':2,'northwest':3}}, inplace=True)

In [24]: insurance_data

Out[24]:
   age  sex  bmi  children  smoker  region  charges
0   19    1  27.900    0      0    1  16884.92400
1   18    0  33.770    1      1    0  1725.55230
2   28    0  33.000    3      1    0  4449.46200
3   33    0  22.705    0      1    3  21984.47061
4   32    0  28.880    0      1    3  3866.85520

...

1333  50    0  30.970    3      1    3  10600.54830
1334  18    1  31.920    0      1    2  2205.98080
1335  18    1  36.850    0      1    0  1629.83350
1336  21    1  25.800    0      1    1  2007.94500
1337  61    1  29.070    0      0    3  29141.36030

1338 rows x 7 columns
```

```
In [40]: # Assuming insurance_data is a Pandas DataFrame
correlation_matrix = insurance_data.corr()

# Create a heatmap of the correlation matrix
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm")

# Set plot labels and title
plt.xlabel('Features')
plt.ylabel('Features')
plt.title('Correlation Heatmap')

# Display the Heatmap
plt.show()

Correlation Heatmap
age      1  0.021  0.11  0.042  0.025  0.0052  0.3
sex      0.021  1 -0.046 -0.017  0.076  0.016 -0.057
bmi      0.11 -0.046  1  0.013 -0.0038 -0.26  0.2
charges  0.042 -0.017  0.013  1 -0.0077  0.019  0.068
children 0.025  0.076 -0.0038 0.0077  1  0.054 -0.79
smoker   0.0052  0.016 -0.26  0.019  0.054  1 -0.057
region   0.3 -0.057  0.2  0.068 -0.79 -0.057  1
Features
age sex bmi children smoker region charges
```

Splitting the Features and Target

```
In [25]: X = insurance_data.drop(columns='charges', axis=1)
Y = insurance_data['charges']

In [26]: print(X)

   age  sex  bmi  children  smoker  region
0   19    1  27.900    0      0    1
1   18    0  33.770    1      1    0
2   28    0  33.000    3      1    0
3   33    0  22.705    0      1    3
4   32    0  28.880    0      1    3
...
1333  50    0  30.970    3      1    3
1334  18    1  31.920    0      1    2
1335  18    1  36.850    0      1    0
1336  21    1  25.800    0      1    1
1337  61    1  29.070    0      0    3

[1338 rows x 6 columns]

In [27]: print(Y)

16884.92400
1725.55230
4449.46200
21984.47061
3866.85520
...
10600.54830
2205.98080
1629.83350
2007.94500
29141.36030
Name: charges, Length: 1338, dtype: float64
```

Splitting the data into Training data & Testing Data

```
In [28]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)

In [29]: print(X.shape, X_train.shape, X_test.shape)

(1338, 6) (1076, 6) (268, 6)
```

Model Training

Linear Regression

```
In [30]: # loading the Linear Regression model
regressor = LinearRegression()

In [31]: regressor.fit(X_train, Y_train)
```

```
Out[31]:
LinearRegression()
```

Model Evaluation

```
In [32]: # prediction on training data
training_data_prediction = regressor.predict(X_train)

In [33]: # R squared value
r2_train = metrics.r2_score(Y_train, training_data_prediction)
print('R squared value : ', r2_train)
R squared value : 0.75159564341174

In [34]: # prediction on test data
test_data_prediction = regressor.predict(X_test)

In [35]: # R squared value
r2_test = metrics.r2_score(Y_test, test_data_prediction)
print('R squared value : ', r2_test)
R squared value : 0.744727389984077

In [45]: plt.scatter(Y_test, test_data_prediction)
plt.xlabel("Actual Charges")
plt.ylabel("Predicted Charges")
plt.show()

Predicted Charges
40000
30000
20000
10000
0
0 10000 20000 30000 40000 50000 60000
Actual Charges
```

Building a Predictive System

```
In [46]: input_data = (31,1,25.74,0,1,0)

# changing input_data to a numpy array
input_data_as_numpy_array = np.asarray(input_data)

# reshape the array
input_data_reshaped = input_data_as_numpy_array.reshape(1,-1)

prediction = regressor.predict(input_data_reshaped)
print(prediction)

[3769.0895765]

The insurance cost is USD ', prediction[0])

C:\Users\hvh\AppData\Local\Temp\ipykernel_13789\base.py:439: UserWarning: X does not have valid feature names, but LinearRegression was fitted with feature names
warnings.warn()
```