

# Medical Cost Personal Insurance Project

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
warnings.filterwarnings('ignore')

In [3]: df=pd.read_csv('medical_cost_insurance.csv')
df.head()

Out[3]:
```

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	1	yes	southwest	16884.92400
1	18	male	33.770	1	no	southwest	1725.55230
2	28	male	33.000	3	no	southeast	4448.40200
3	31	male	22.765	0	no	northwest	21984.47060
4	32	male	28.880	0	no	northwest	3968.85520

```
In [4]: df.tail()

Out[4]:
```

	age	sex	bmi	children	smoker	region	charges
1333	50	male	30.97	3	no	northwest	10600.5483
1334	18	female	31.92	0	no	northeast	2205.9088
1335	11	female	24.85	0	no	southeast	1829.8385
1336	21	female	25.60	0	no	southeast	2927.8460
1337	61	female	29.07	0	yes	northwest	29141.3503

```
In [5]: df.shape
(1338, 7)

In [6]: df.describe()

Out[6]:
```

	age	bmi	children	charges
count	1338.000000	1338.000000	1338.000000	1338.000000
mean	39.207025	30.663397	1.094818	13270.422265
std	14.049900	6.096187	1.254493	12110.011237
min	18.000000	15.960000	0.000000	1123.873900
25%	27.000000	26.296250	0.000000	4740.287150
50%	39.000000	30.663397	1.000000	16884.924000
75%	46.000000	34.499250	2.000000	16629.92515
max	64.000000	53.130000	5.000000	63770.43010

```
In [7]: df.dtypes

Out[7]:
```

	age	sex	bmi	children	smoker	region	charges
age	int64						
sex	object						
bmi	float64						
children	int64						
smoker	object						
region	object						
charges	float64						
dtype:	object						

```
In [8]: df.isnull().sum()

Out[8]:
```

	age	sex	bmi	children	smoker	region	charges
age	0						
sex	0						
bmi	0						
children	0						
smoker	0						
region	0						
charges	0						
dtype:	int64						

## EDA and Visualizations

```
In [9]: sns.set(style='whitegrid')
f, ax = plt.subplots(1,1,figsize=(12, 8))
ax = sns.distplot(df['charges'], kde = True, color = 'c')
plt.title('Distribution of Charges')

Out[9]:
```

```
In [10]: f, ax = plt.subplots(1, 1, figsize=(12, 8))
ax = sns.distplot(np.log(df['charges']), kde = True, color = 'r')

Out[10]:
```

```
In [11]: #How we are calculating total charges by region column
#Sorting it by ascending value and creating a bar plot to visualize the top few regions with the lowest charges
charges = df.groupby('region')['charges'].sort_values(ascending=True)
f, ax = plt.subplots(1, 1, figsize=(8,6))
ax = sns.barplot(x=charges.index, y=charges.head(), palette='Blues')

Out[11]:
```

```
In [12]: #How I want to compare the charges across different regions with the bars further grouped by sex using hue parameter
f, ax = sns.barplot(x='region', y='charges', hue='sex', data=df, palette='cool')

Out[12]:
```

```
In [13]: # How I want to create a bar plot to compare the charges across different regions with the bars grouped by smoker using the hue parameter
f, ax = sns.barplot(x='region', y='charges', hue='smoker', data=df, palette='Reds_r')

Out[13]:
```

```
In [14]: #How I am creating a bar plot to compare the charges across different regions with the bars further grouped by the number of children using the hue
f, ax = plt.subplots(1, 1, figsize=(12, 8))
ax = sns.barplot(x='region', y='charges', hue='children', data=df, palette='Set1')

Out[14]:
```

```
In [15]: #How I am creating three separate scatter plots to visualize the relationship between age bmi children with charges
ax = sns.lmplot(x='age', y='charges', data=df, hue='smoker', palette='Set1')
ax = sns.lmplot(x='bmi', y='charges', data=df, hue='smoker', palette='Set2')
ax = sns.lmplot(x='children', y='charges', data=df, hue='smoker', palette='Set3')

Out[15]:
```

```
In [16]: #How I am using violinplot to visualize the distribution of charges on the number of children
f, ax = plt.subplots(1,1,figsize=(10, 10))
ax = sns.violinplot(x='children', y='charges', data=df, orient='v', hue='smoker', palette='inferno')

Out[16]:
```

```
In [17]: #Converting objects labels into categorical
df[['sex', 'smoker', 'region']] = df[['sex', 'smoker', 'region']].astype('category')
df.dtypes

Out[17]:
```

	age	sex	bmi	children	smoker	region	charges
age	int64						
sex	category						
bmi	float64						
children	int64						
smoker	category						
region	category						
charges	float64						
dtype:	object						

```
In [18]: #Converting category labels into numerical using LabelEncoder
from sklearn.preprocessing import LabelEncoder
label = LabelEncoder()
label.fit(df.sex.drop_duplicates())
df.sex = label.transform(df.sex)
label.fit(df.smoker.drop_duplicates())
df.smoker = label.transform(df.smoker)
label.fit(df.region.drop_duplicates())
df.region = label.transform(df.region)
df.dtypes

Out[18]:
```

	age	sex	bmi	children	smoker	region	charges
age	int64						
sex	int32						
bmi	float64						
children	int64						
smoker	int32						
region	int64						
charges	float64						
dtype:	object						

```
In [19]: f, ax = plt.subplots(1, 1, figsize=(10, 10))
ax = sns.heatmap(df.corr(), annot=True, cmap='cool')

Out[19]:
```

## Training model

### Linear Regression

```
In [20]: #Importing the necessary libraries
from sklearn.model_selection import train_test_split as holdout
from sklearn.linear_model import LinearRegression
from sklearn import metrics

#Load the dataset
x = df.drop(['charges'],axis=1)
y = df['charges']

#Split the dataset into training and testing datasets
x_train,x_test,y_train,y_test = holdout(x,y, test_size=0.2, random_state=0)

#Create the Linear Regression model
lin_reg = LinearRegression()
lin_reg.fit(x_train,y_train)

#Print the intercept,coeff, and score of the model
print(lin_reg.intercept_)
print(lin_reg.coef_)
print(lin_reg.score(x_test,y_test))

-1163.88398882442
[ 2.3395148e+02  3.3455988e+01  3.2833956e+02  4.4422847e+02
  2.3869957e+04 -2.8865857e+02]
0.88129892383228
```

### Ridge Regression

```
In [21]: #Importing the required library
from sklearn.linear_model import Ridge

#Creating the Ridge Regression Model
Ridge = Ridge(alpha=0.5)

#Train the Ridge Regression Model
Ridge.fit(x_train,y_train)

#Print the intercept,coeff, and score of the model
print(Ridge.intercept_)
print(Ridge.coef_)
print(Ridge.score(x_test,y_test))

-1163.449927495936
[ 2.3395148e+02  3.3455988e+01  3.2833956e+02  4.4422847e+02
  2.3869957e+04 -2.8865857e+02]
0.88129892383228
```

### Lasso Regression

```
In [22]: #Importing the library
from sklearn.linear_model import Lasso

#Creating the Lasso Regression Model
Lasso = Lasso(alpha=0.2, fit_intercept=True, precompute=False, max_iter=1000,
              tol=0.0001, warm_start=False, positive=False, random_state=None, selection='cyclic')

#Train the Lasso Regression Model
Lasso.fit(x_train,y_train)

#Print the intercept,coeff, and score of the model
print(Lasso.intercept_)
print(Lasso.coef_)
print(Lasso.score(x_test,y_test))

-1163.88398882442
[ 2.3395148e+02  3.3455988e+01  3.2833956e+02  4.4422847e+02
  2.3869957e+04 -2.8865857e+02]
0.88129892383228
```

## Random Forest Regressor

```
In [23]: #Importing the library
from sklearn.ensemble import RandomForestRegressor as rfr

#Load the dataset
x = df.drop(['charges'],axis=1)
y = df['charges']

#Create a Random Forest Regression
Rfr = rfr(n_estimators=100, criterion='squared_error',
         random_state=None, n_jobs=-1)

#Train the Random Forest Model
Rfr.fit(x_train,y_train)

#Make predictions of the model
x_train_pred = Rfr.predict(x_train)
x_test_pred = Rfr.predict(x_test)

#Print the mean square of the model
print('MSE train data: %.3f, MSE test data: %.3f' %
      (metrics.mean_squared_error(x_train_pred,y_train),
       metrics.mean_squared_error(x_test_pred,y_test)))

MSE train data: 3628893.887, MSE test data: 19730453.918
```

## Results and Analysis

```
In [24]: plt.figure(figsize=(8,6))

plt.scatter(x_train_pred,x_train_pred - y_train,
            c='gray', marker='o', s=35, alpha=0.5,
            label='Train data')
plt.scatter(x_test_pred,x_test_pred - y_test,
            c='blue', marker='o', s=35, alpha=0.7,
            label='Test data')

plt.xlabel('Predicted values')
plt.ylabel('Actual values')
plt.legend(loc='upper right')
plt.xlim(y.min()*0.5, y.max()*1.5)
plt.ylim(-20000, 20000)

#matplotlib.collections.LineCollection at 0x246e4574a80

Out[24]:
```

```
In [33]: print('Feature Importance Ranking\n')
importances = Rfr.feature_importances_
std = np.std([tree.feature_importances_ for tree in Rfr.estimators_],axis=0)
indices = np.argsort(importances)[::-1]
for f in range(x.shape[0]):
    print("%5s (%7s)" % (f+1, variable, importances[indices[f]]))
    variable = variables[indices[f]]
    importance_list.append(variable)
    print("\n\n")

#Plot the Feature importance of the forest
plt.figure()
plt.title('Feature Importances')
plt.bar(importance_list, importances[indices],
        color='y', yerr=std[indices], align='center')

Feature importance ranking

1.smoker(0.881678)
2.bmi(0.218831)
3.age(0.139943)
4.children(0.880999)
5.region(0.815165)
6.smoker(0.88093)
<BarContainer object of 6 artists>
```

```
In [34]: #Importing the library
from sklearn.preprocessing import PolynomialFeatures

#Load the data set
x = df.drop(['sex','smoker','region'],axis=1)
y = df['charges']

#Create polynomial features
pol = PolynomialFeatures(degree=2)
x_pol = pol.fit_transform(x)

#Split the dataset into training and testing
x_train,x_test,y_train,y_test = holdout(x_pol,y, test_size=0.2, random_state=0)

#Create the polynomial regression model
pol_reg = LinearRegression()
pol_reg.fit(x_train,y_train)

#Make predictions of the model
y_train_pred = pol_reg.predict(x_train)
y_test_pred = pol_reg.predict(x_test)

#Print the intercept,coeff, and score of the model
print(pol_reg.intercept_)
print(pol_reg.coef_)
print(pol_reg.score(x_test,y_test))

-1163.88398882442
[ 2.3395148e+02  3.3455988e+01  3.2833956e+02  4.4422847e+02
  2.3869957e+04 -2.8865857e+02]
0.88129892383228
```

```
In [37]: #Evaluating the charges
y_train_pred = pol_reg.predict(x_test)
#Comparing the actual output values with the predicted values
df = pd.DataFrame({'Actual': y_test, 'Predicted': y_test_pred})

Out[37]:
```

	Actual	Predicted
578	8747.613000	12101.150236
619	47702.02235	40541.022961
1354	12960.07120	34163.007562
136	8644.25229	8636.230727
...	...	...
1084	15019.76055	16712.196261
728	6664.68695	8654.956461
1132	20709.02034	12372.050009
729	40932.42952	41465.617288
963	9500.97385	10941.780765

```
In [38]: #Importing the library
from sklearn.preprocessing import PolynomialFeatures

#Load the data set
x = df.drop(['sex','smoker','region'],axis=1)
y = df['charges']

#Create polynomial features
pol = PolynomialFeatures(degree=2)
x_pol = pol.fit_transform(x)

#Split the dataset into training and testing
x_train,x_test,y_train,y_test = holdout(x_pol,y, test_size=0.2, random_state=0)

#Create the polynomial regression model
pol_reg = LinearRegression()
pol_reg.fit(x_train,y_train)

#Make predictions of the model
y_train_pred = pol_reg.predict(x_train)
y_test_pred = pol_reg.predict(x_test)

#Print the intercept,coeff, and score of the model
print(pol_reg.intercept_)
print(pol_reg.coef_)
print(pol_reg.score(x_test,y_test))

-1163.88398882442
[ 2.3395148e+02  3.3455988e+01  3.2833956e+02  4.4422847e+02
  2.3869957e+04 -2.8865857e+02]
0.88129892383228
```

```
In [39]: #Evaluating the charges
y_train_pred = pol_reg.predict(x_test)
#Comparing the actual output values with the predicted values
df = pd.DataFrame({'Actual': y_test, 'Predicted': y_test_pred})

Out[39]:
```

	Actual	Predicted
578	8747.613000	12101.150236
619	47702.02235	40541.022961
1354	12960.07120	34163.007562
136	8644.25229	8636.230727
...	...	...
1084	15019.76055	16712.196261
728	6664.68695	8654.956461
1132	20709.02034	12372.050009
729	40932.42952	41465.617288
963	9500.97385	10941.780765

```
In [39]: #Evaluating the charges
y_train_pred = pol_reg.predict(x_test)
#Comparing the actual output values with the predicted values
df = pd.DataFrame({'Actual': y_test, 'Predicted': y_test_pred})

Out[39]:
```

	Actual	Predicted
578	8747.613000	12101.150236
619	47702.02235	40541.022961
1354	12960.07120	34163.007562
136	8644.25229	8636.230727
...	...	...
1084	15019.76055	16712.196261
728	6664.68695	8654.956461
1132	20709.02034	12372.050009
729	40932.42952	41465.617288
963	9500.97385	10941.780765

## Polynomial Regression

```
In [38]: #Importing the library
from sklearn.preprocessing import PolynomialFeatures

#Load the data set
x = df.drop(['sex','smoker','region'],axis=1)
y = df['charges']

#Create polynomial features
pol = PolynomialFeatures(degree=2)
x_pol = pol.fit_transform(x)

#Split the dataset into training and testing
x_train,x_test,y_train,y_test = holdout(x_pol,y, test_size=0.2, random_state=0)

#Create the polynomial regression model
pol_reg = LinearRegression()
pol_reg.fit(x_train,y_train)

#Make predictions of the model
y_train_pred = pol_reg.predict(x_train)
y_test_pred = pol_reg.predict(x_test)

#Print the intercept,coeff, and score of the model
print(pol_reg.intercept_)
print(pol_reg.coef_)
print(pol_reg.score(x_test,y_test))

-1163.88398882442
[ 2.3395148e+02  3.3455988e+01  3.2833956e+02  4.4422847e+02
  2.3869957e+04 -2.8865857e+02]
0.88129892383228
```

```
In [37]: #Evaluating the charges
y_train_pred = pol_reg.predict(x_test)
#Comparing the actual output values with the predicted values
df = pd.DataFrame({'Actual': y_test, 'Predicted': y_test_pred})

Out[37]:
```

	Actual	Predicted
578	8747.613000	12101.150236
619	47702.02235	40541.022961
1354	12960.07120	34163.007562
136	8644.25229	8636.230727
...	...	...
1084	15019.76055	16712.196261
728	6664.68695	8654.956461
1132	20709.02034	12372.050009
729	40932.42952	41465.617288
963	9500.97385	10941.780765

```
In [38]: #Importing the library
from sklearn.preprocessing import PolynomialFeatures

#Load the data set
x = df.drop(['sex','smoker','region'],axis=1)
y = df['charges']

#Create polynomial features
pol = PolynomialFeatures(degree=2)
x_pol = pol.fit_transform(x)

#Split the dataset into training and testing
x_train,x_test,y_train,y_test = holdout(x_pol,y, test_size=0.2, random_state=0)

#Create the polynomial regression model
pol_reg = LinearRegression()
pol_reg.fit(x_train,y_train)

#Make predictions of the model
y_train_pred = pol_reg.predict(x_train)
y_test_pred = pol_reg.predict(x_test)

#Print the intercept,coeff, and score of the model
print(pol_reg.intercept_)
print(pol_reg.coef_)
print(pol_reg.score(x_test,y_test))

-1163.88398882442
[ 2.3395148e+02  3.3455988e+01  3.2833956e+02  4.4422847e+02
  2.3869957e+04 -2.8865857e+02]
0.88129892383228
```

```
In [37]: #Evaluating the charges
y_train_pred = pol_reg.predict(x_test)
#Comparing the actual output values with the predicted values
df = pd.DataFrame({'Actual': y_test, 'Predicted': y_test_pred})

Out[37]:
```

	Actual	Predicted
578	8747.613000	12101.150236
619	47702.02235	40541.022961
1354	12960.07120	34163.007562
136	8644.25229	8636.230727
...	...	...
1084	15019.76055	16712.196261
728	6664.68695	8654.956461
1132	20709.02034	12372.050009
729	40932.42952	41465.617288
963	9500.97385	10941.780765

```
In [38]: #Importing the library
from sklearn.preprocessing import PolynomialFeatures

#Load the data set
x = df.drop(['sex','smoker','region'],axis=1)
y = df['charges']

#Create polynomial features
pol = PolynomialFeatures(degree=2)
x_pol = pol.fit_transform(x)

#Split the dataset into training and testing
x_train,x_test,y_train,y_test = holdout(x_pol,y, test_size=0.2, random_state=0)

#Create the polynomial regression model
pol_reg = LinearRegression()
pol_reg.fit(x_train,y_train)

#Make predictions of the model
y_train_pred = pol_reg.predict(x_train)
y_test_pred = pol_reg.predict(x_test)

#Print the intercept,coeff, and score of the model
print(pol_reg.intercept_)
print(pol_reg.coef_)
print(pol_reg.score(x_test,y_test))

-1163.88398882442
[ 2.3395148e+02  3.3455988e+01  3.2833956e+02  4.4422847e+02
  2.3869957e+04 -2.8865857e+02]
0.88129892383228
```

```
In [37]: #Evaluating the charges
y_train_pred = pol_reg.predict(x_test)
#Comparing the actual output values with the predicted values
df = pd.DataFrame({'Actual': y_test, 'Predicted': y_test_pred})

Out[37]:
```

	Actual	Predicted
578	8747.613000	12101.150236
619	47702.02235	40541.022961
1354	12960.07120	34163.007562
136	8644.25229	8636.230727
...	...	...
1084	15019.76055	16712.196261
728	6664.68695	8654.956461
1132	20709.02034	12372.050009
729	40932.42952	41465.617288
963	9500.97	