

```
[In [1]]: # MEDICAL INSURANCE DATASET
# Is the Machine Learning task for this dataset is classification as it has labelled data which enables ML algorithm to find the relationship between any two points(variables)

[In [2]]: # i.e. The Machine Learning task here for this dataset is classification as it has labelled data. Logistic regression is not the best type of classification for this dataset as it is best for making a prediction about a categorical variable.

[In [3]]: #ii. Data Exploration : When observing the dataset, it can be seen that the older the individual, the higher the medical cost. Also, the higher the bmi, the higher the medical cost of those people. However, explanatory data tools are required to confirm this.

[In [4]]: # To start this assessment, data exploration comes first using the relevant data exploration tools, some relevant libraries would be imported and the dataset uploaded

[In [5]]:
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
style.use('seaborn-whitegrid')
plt.rcParams['figure.figsize'] = (20,10)
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score as r2
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score

[In [6]]:
dataset = pd.read_csv('insurance.csv')
dataset

Out[6]:
   age  sex  bmi  children  smoker  region  medicalCost
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
...
1332 61  female  29.070    0    yes  northwest  29141.36310
1338 rows x 7 columns

[In [7]]:
# Exploration of the dataset
dataset.isnull()

Out[7]:
   age  sex  bmi  children  smoker  region  medicalCost
0  False  False  False  False  False  False  False
1  False  False  False  False  False  False  False
2  False  False  False  False  False  False  False
3  False  False  False  False  False  False  False
4  False  False  False  False  False  False  False
...
1332  False  False  False  False  False  False  False
1338 rows x 7 columns

[In [8]]:
# No null values. The dataset has 7 Columns and 1338 rows

[In [9]]:
# Check the type of data whether integers or float
dataset.dtypes

Out[9]:
age      int64
sex      object
bmi      float64
children  int64
smoker    object
region    object
medicalCost  float64
dtype: object

[In [10]]:
# Only bmi and smoker cat columns have float variables, the rest have objects

[In [11]]:
# First obtain the statistical values of all the variables

[In [12]]:
dataset.describe()

Out[12]:
   age      bmi  children  medicalCost
count  1338.000000  1338.000000  1338.000000  1338.000000
mean     29.207025   30.663397   1.094818  12705.422265
std     14.649660    6.998187   1.205491  12110.012137
min      18.000000   15.960000    0.000000   1121.873900
25%     27.000000   26.296150    0.000000   4740.287150
50%     29.000000   30.400000   1.000000   9382.033000
75%     31.000000   34.693750   2.000000  16639.912515
max     64.000000   51.130000   5.000000  63776.426010

[In [13]]:
# mean value for age column is 29, bmi mean is 30.66, children column mean value is 1.09 and medical cost mean value is 12705

[In [14]]:
# check for null values

[In [15]]:
dataset.isnull()

Out[15]:
   age  sex  bmi  children  smoker  region  medicalCost
0  False  False  False  False  False  False  False
1  False  False  False  False  False  False  False
2  False  False  False  False  False  False  False
3  False  False  False  False  False  False  False
4  False  False  False  False  False  False  False
...
1332  False  False  False  False  False  False  False
1338 rows x 7 columns

[In [16]]:
# view the first 5 rows of the dataset

[In [17]]:
dataset.head()

Out[17]:
   age  sex  bmi  children  smoker  region  medicalCost
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

[In [18]]:
# Logistic regression as a machine model will not be best for this dataset as it can work best only with categorical variables, as seen below:
# dataset = {'sex': 'female', 'region': 'southwest', 'medicalCost': 'age', 'bmi': 1, axis = 1}
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
from sklearn.linear_model import LogisticRegression
lr_model = LogisticRegression()
lr_model.fit(X_train, y_train)

ValueError                                Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_1664\321641604.py in <module>
      5 from sklearn.linear_model import LogisticRegression
----> 6 lr_model = LogisticRegression()
      7 lr_model.fit(X_train, y_train)
----> 7 lr_model.fit(X_train, y_train)

~\anaconda3\lib\site-packages\sklearn\linear_model\logistic.py in fit(self, X, y, sample_weight)
   1345         order='C')
   1346         accept_sparse = sp.sparse.spmatrix
--> 1347         check_classification_targets(y)
   1348         self.classes_ = np.unique(y)
   1349

~\anaconda3\lib\site-packages\sklearn\utils\multiclass.py in check_classification_targets(y)
   182         'multilabel-indicator', 'multilabel-sequences'])
--> 183         raise ValueError('Unknown label type: %s' % y.dtype)
   184
   185

ValueError: Unknown label type: 'continuous'
```