Problem Statement or Requirement:

A client's requirement is, he wants to predict the insurance charges based on the several parameters. The Client has provided the dataset of the same.

Data set Analysis:

1. Domain selection: Machine Learning

Learning selection: Supervised Learning

Method: Regression

- 2. Dataset contains 6 columns: Age, sex, BMI, children, smoker, charge &1138 rows
- 3. I could see 2 nominal columns: sex and smoker which must be converted with the help of one hot encoding

Since the data involves multiple inputs, here we can't use Simple Linear Regression. So, let's start with other available algorithms.

1.MLR - Multiple Linear Regression

R2_score value: 0.7894790349867009

```
14410.20040
     410
     920
          13451.12200
     [402 rows x 1 columns],
     (402, 1))
from sklearn.linear_model import LinearRegression
    regressor=LinearRegression()
    regressor.fit(x_train,y_train)
]:

    LinearRegression

   LinearRegression()
]: weight=regressor.coef_
    weight
                           321.06004271,
]: array([[ 257.8006705 ,
                                             469.58113407,
                                                            -41.74825718,
           23418.6671912 ]])
    bais=regressor.intercept_
    bais
]: array([-12057.244846])
  y_pred=regressor.predict(x_test)
]: from sklearn.metrics import r2_score
    R_value=r2_score(y_test,y_pred)
    R_value
0.7894790349867009
```

2. Support Vector Machine – SVM:

R2 score value (without parameter tuning): -0.08338238593619329

With Parameter turning:

S.NO	Hyper	Linear	Non-Linear	Poly	Sigmoid
	parameter		(RBF)		
1	C=0.01	- 0.088831334391 68489	- 0.08964553739867 864	- 0.089568284876 71076	- 0.089565015934 1983
2	C=0.10	- 0.080959968427 891	- 0.08907451521042 731	- 0.088302376554 10711	- 0.088269914504 85111
3	C=1	- 0.01010266531608 1394	- 0.08338238593619 329	- 0.075699655708 60893	- 0.075429242811 07188
4	C=10	0.46246841423396 834	- 0.03227329390671 052	0.038716222760 231456	0.039307143782 74347
5	C=100	0.62887928573203 69	0.32003178320508 31	0.617956962405 9795	0.527610354651 0407
6	C=1000	0.76493117385974 11	0.81020648517585 45	0.856648767594 6572	0.287470694869 76173
7	C=2000	0.74404183081078 46	0.85477664253929 79	0.860557925859 7704	- 0.593950973128 3505

```
[402 rows x 1 columns],
        (402, 1))
[217]: from sklearn.svm import SVR
       regressor=SVR(kernel='poly',C=2000)
       regressor.fit(x_train,y_train)
       C:\Anaconda\Lib\site-packages\sklearn\utils\validation.py:1408: DataConversionWarning: A co
       lumn-vector y was passed when a 1d array was expected. Please change the shape of y to (n_s)
       amples, ), for example using ravel().
        y = column_or_1d(y, warn=True)
                 SVR
       SVR(C=2000, kernel='poly')
[218]: regressor.intercept_
[218]: array([8281.60615314])
[219]: regressor.n_support_
[219]: array([936], dtype=int32)
[220]: y_pred=regressor.predict(x_test)
[221]: from sklearn.metrics import r2_score
       R_value=r2_score(y_test,y_pred)
[222]: R_value
[222]: 0.8605579258597704
```

3. Decision Tree:

R2 score value (without parameter tuning): 0.6966865293326343

With Parameter turning:

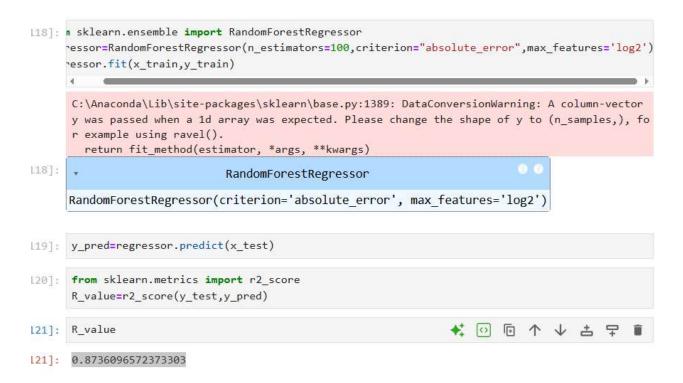
s.no	criterion	Max feature	splitter	R score
1.	Squared error	sqrt	best	0.702035228596572
2.	Squared error	log2	best	0.7419054872546678
3.	Squared error	sqrt	random	0.682074724492034
4.	Squared error	log2	random	0.6980896260139877
5.	Squared error	none	best	0.6996994481671908
6.	Squared error	none	random	0.7358200878574843
7.	Friedman mse	sqrt	best	0.7381415595482762
8.	Friedman mse	log2	best	0.7319704870186678
9.	Friedman mse	sqrt	random	0.6457245773642908
10.	Friedman mse	log2	random	0.688878706147374
11.	Friedman mse	none	best	0.6879110478137831
12.	Friedman mse	none	random	0.6949890484996232
13.	MAE	sqrt	best	0.7297458122874192
14.	MAE	log2	best	0.6420218622840612
<mark>15.</mark>	MAE	<mark>sqrt</mark>	<mark>random</mark>	<mark>0.7607657594740639</mark>
16.	MAE	log2	random	0.7240231565275986
17.	MAE	none	best	0.6563984953703915
18.	MAE	none	random	0.6888543930260137
19.	Poisson	sqrt	best	0.7052009320217727
20.	Poisson	log2	best	0.6199992359174482
21.	Poisson	sqrt	random	0.646019598776802
22.	Poisson	log2	random	0.6597185842121801
23.	Poisson	none	best	0.7240970726548698
24.	Poisson	none	random	0.7360710221488609

4. Random Forest:

R2 score value (without parameter tuning): 0.857238919537012

With Parameter turning:

s.no	N estimators	criterion	Max features	R score
1	10	Squared error		0.8297825655643111
2	100	Squared error		0.8575979556497877
3	10	Friedman mse		0.8318936415190324
4	100	Friedman mse		0.8480668109283795
5	10	MAE		0.8368339827516286
6	100	MAE		0.856093154513048
7	10	Poisson		0.8274737517834232
8	100	Poisson		0.8540244751233814
9	100	Squared error	sqrt	0.8724894738130522
10	100	Squared error	log2	0.8676627351539128
11	100	MAE	sqrt	0.8704424718543277
12	<mark>100</mark>	MAE	log2	0.8736096572373303



Conclusion:

For this dataset, I use different types of algorithm and hyper parameters from this I find the better performing algorithm in **Random Forest.** In that I created a finial model.

The best predicted value of **R_Score** is 0.8736096572373303 in random forest (n_estimators=100, criterion=MAE,max feature=log2).