

## 1. Problem Statement :

Based on the Parameters like Age, Sex, Children,Bmi, Smoker . We need to predict the Insurance charge for the Person.

Stage 1 - Machine Learning (as data deals with numbers )

Stage 2 - Supervised learning (requirement and inputs/outputs are clear )

Stage 3 - Regression (final output is in numerical data)

## 2. Basic Information about Dataset :

The dataset has 1338 rows and 6 columns

The Input variables are Age, Sex, Bmi, Children, Smoker

The Ouput Variable is Insurance Charges

## 3. Pre Processing method used :

In the dataset the columns Smoker and Sex are categorical and it is Nominal Type .  
So **One Hot Encoding** is used to convert it into numerical

## R2 Scores of the Models :

1. Multiple linear Regression : **0.789479**

## 2. Support Vector Machine

Hyper Parameter	linear	poly	rbf	sigmoid
C = 10	0.5665127	0.159392	-0.01810	0.073055
C = 100	0.6359503	0.750819	0.390601	0.52756
C = 500	0.76514321	0.859312	0.696468	0.490635
C = 1000	0.7440909	0.860584	0.828352	0.14377
C = 2000	0.74142388	0.860181	0.860735	-2.58403
C = 3000	0.74142316	0.860011	0.868531	-6.82618

SVM : For 'rbf' and C = 3000 the model gives **0.868531** as R score

## 3. Decision Tree :

critierion	max_features	splitter	R_Score
<i>squared_error</i>	-	Best	0.7129
<i>friedman_mse</i>	-	Best	0.6909
<i>absolute_error</i>	-	Best	0.6796
<i>poisson</i>	-	Best	0.7258
<i>squared_error</i>	-	random	0.7189

<i>friedman_mse</i>	-	random	0.7431
<i>absolute_error</i>	-	random	0.7551
<i>poisson</i>	-	random	0.6521
<i>squared_error</i>	sqrt	Best	0.7297
<i>friedman_mse</i>	sqrt	Best	0.7332
<i>absolute_error</i>	sqrt	Best	0.7434
<i>poisson</i>	sqrt	Best	0.7150
<i>squared_error</i>	sqrt	random	0.7420
<i>friedman_mse</i>	sqrt	random	0.7413
<i>absolute_error</i>	sqrt	random	0.7178
<i>poisson</i>	sqrt	random	0.6535
<i>squared_error</i>	log2	Best	0.7082
<i>friedman_mse</i>	log2	Best	0.7427
<i>absolute_error</i>	log2	Best	0.6822
<i>poisson</i>	log2	Best	0.7634
<i>squared_error</i>	log2	random	0.6907
<i>friedman_mse</i>	log2	random	0.6582
<i>absolute_error</i>	log2	random	0.6826
<i>poisson</i>	log2	random	0.6194

The Good R score is obtained for (Poisson,Log2,best) = **0.7634**

#### 4. Random Forest

<b>criterion</b>	<b>max_features</b>	<b>n_estimators</b>	<b>R_Score</b>
<i>squared_error</i>	<i>log2</i>	50	0.8667
<i>friedman_mse</i>	<i>log2</i>	50	0.8653
<i>absolute_error</i>	<i>log2</i>	50	0.8654
<i>poisson</i>	<i>log2</i>	50	0.8670
<i>squared_error</i>	<i>log2</i>	100	0.8698
<i>friedman_mse</i>	<i>log2</i>	100	0.8681
<i>absolute_error</i>	<i>log2</i>	100	0.8657
<i>poisson</i>	<i>log2</i>	100	0.8687
<i>squared_error</i>	<i>sqrt</i>	50	0.8667
<i>friedman_mse</i>	<i>sqrt</i>	50	0.8653
<i>absolute_error</i>	<i>sqrt</i>	50	0.8653
<i>poisson</i>	<i>sqrt</i>	50	0.8670
<i>squared_error</i>	<i>sqrt</i>	100	0.8698
<i>friedman_mse</i>	<i>sqrt</i>	100	0.8681
<i>absolute_error</i>	<i>sqrt</i>	100	0.8657
<i>poisson</i>	<i>sqrt</i>	100	0.8687

RF The Model performed well in (Squared\_error, sqrt, 100) = **0.8698**

**The final model that can be chosen is**

**SVM ('rbf' , C = 3000) = 0.868531 or**

**RF (Squared\_error, sqrt, 100) = 0.8698**

