

### 1.) Identify the problem statement

- Dataset is very clear and provided input and output data,so this is supervised learning
- We need to predict the insurance charges,so this will come under regression
- There are so many input so we will not use simple linear regression

### 2.) About the dataset

- Dataset have 1337 rows and 6 column

### 3.) pre-processing the dataset

- There are 2 string columns that need to be converting into numeric
- Columns Sex and smoker are non orderable so it will come under nominal data
- For nominal data we need to use one-hot coding method
- Input columns are 'age', 'bmi', 'children', 'sex\_male', 'smoker\_yes'
- Output column is 'charges'

### R2 Value:

- Multiple linear regression we got R2 value is 0.78
- SVM Module R2 value tuning using below parameter

## Finding Best SVM Module using R2 value

kernel	C	coef0	degree	R2 value
<i>linear</i>	1000			0.76
<i>poly</i>	1000			0.85
<i>rbf</i>	1000			0.21
<i>sigmoid</i>	1000			0.85
<i>sigmoid</i>	1000	0.5		-1.62
<i>poly</i>	1000	0.5		0.87
<i>poly</i>	1000	0.5	3	0.87
<i>poly</i>	1000	0.5	1	0.63
<i>poly</i>	1000		1	0.87
<i>poly</i>	1000	155.5		0.86
<i>linear</i>	2000			0.74

We got high r2 value = 0.87 for kernel=poly ,C=1000 parameter and coef0=0.5,so this will be the best model

## Finding Best Decision Tree Module using R2 value

R2 will change depends up on various hyper tuning parameter

critierion	splitter	min_samples_split	ccp_alpha	random_state	R2 value
<i>squared_error</i>					0.67
<i>squared_error</i>	best	2	0.0	None	0.67
<i>squared_error</i>	<i>random</i>				0.7
<i>squared_error</i>	<i>random</i>	14	0.5		0.84
<i>squared_error</i>	<i>random</i>	5	0.5	10	0.85
<i>friedman_mse</i>					0.69
<i>friedman_mse</i>	best	2	0.0	None	0.69
<i>friedman_mse</i>	<i>random</i>				0.73
<i>friedman_mse</i>	<i>random</i>	14			0.84
<i>friedman_mse</i>	<i>random</i>	14	0.5		0.85
<i>friedman_mse</i>	<i>random</i>	14	0.5	10	0.85
<i>absolute_error</i>					0.69
<i>absolute_error</i>	best	2	0.0	None	0.69
<i>absolute_error</i>	<i>random</i>				0.73

<i>absolute_error</i>	<i>random</i>	14	0.5	10	0.85
<i>absolute_error</i>	<i>random</i>	14			0.87
<i>absolute_error</i>	<i>random</i>	14	0.5		0.88
<i>poisson</i>					0.67
<i>poisson</i>	best	2	0.0	None	0.67
<i>poisson</i>	<i>random</i>				0.71
<i>poisson</i>	<i>random</i>	14			0.82
<i>poisson</i>	<i>random</i>	14	0.5		0.86
<i>poisson</i>	<i>random</i>	14	0.5	10	0.85

We got high r2 value = 0.88 for the hyper tuning parameter  
criterion=absolute\_error,splitter  
=random,min\_samples\_split=14,ccp\_alpha=0.5

## Finding Best Random forest Module using R2 value

R2 will change depends up on various hyper tuning parameter

n_estimators	criterion	random_state	R2 value
	<i>squared_error</i>		0.85
50	<i>squared_error</i>		0.84
50	<i>squared_error</i>	None	0.84
61	<i>squared_error</i>	18	0.84
	<i>friedman_mse</i>		0.85
50	<i>friedman_mse</i>		0.85
61	<i>friedman_mse</i>	None	0.85
61	<i>friedman_mse</i>	18	0.85
	<i>absolute_error</i>		0.85
50	<i>absolute_error</i>		0.85
50	<i>absolute_error</i>	18	0.85
61	<i>absolute_error</i>	18	0.85
	<i>poisson</i>		0.84
50	<i>poisson</i>		0.84
61	<i>poisson</i>	None	0.84

50	<i>poisson</i>	18	0.84
----	----------------	----	------

We got high r2 value =0.85 for the hyper tuning parameter  
criterion=absolute\_error ,friedman\_mse, squared\_error

4. I have created many models using machine learning algorithm like

- Multiple linear regression
- Support vector machine
- Decision tree
- Random forest

We can select the best model having a high R2 value.

Compared to all other values, the Decision Tree R2 value 0.88 is high,so we can finalize that.