

1) Problem Statement

Develop a machine learning model to accurately predict medical insurance charges for new customers based on inputs: age, sex, BMI, children, smoker status, and region.

2) Dataset Information

The dataset contains 1338 rows and 7 columns including both numerical and categorical data.

3) Data Preprocessing

Categorical variables (sex, smoker, region) were encoded numerically. The data was split into training and testing sets, and the target variable was reshaped.

Support Vector Regression Results

S.NO	HYPER PARAMETER	LINEAR	RBF	POLY	SIGMOID
1	C10	0.4625	-0.0323	0.0387	0.0393
2	C100	0.6289	0.32	0.618	0.5276
3	C500	0.7631	0.6643	0.8264	0.4446
4	C1000	0.7649	0.8102	0.8566	0.2875
5	C2000	0.744	0.8548	0.8606	-0.594
6	C3000	0.7414	0.8663	0.8599	-2.1244

Decision Tree Results

SL.NO	CRITERION	MAX FEATURES	SPLITTER	R VALUE
1	squared_error	sqrt	best	0.7082
2	squared_error	sqrt	random	0.7305
3	squared_error	log2	best	0.6945
4	squared_error	log2	random	0.5486
5	absolute_error	sqrt	best	0.7238
6	absolute_error	sqrt	random	0.7603
7	absolute_error	log2	best	0.6398

8	absolute_error	log2	random	0.7225
9	friedman_mse	sqrt	best	0.7036
10	friedman_mse	sqrt	random	0.7242
11	friedman_mse	log2	best	0.5685
12	friedman_mse	log2	random	0.6848

Random Forest Results

SL.NO	CRITERION	MAX FEATURES	N_ESTIMATORS	R VALUE
1	squared_error	Sqrt	10	0.859
2	squared_error	Sqrt	100	0.8706
3	squared_error	Log2	10	0.8484
4	squared_error	Log2	100	0.8697
5	absolute_error	Sqrt	10	0.8447
6	absolute_error	Sqrt	100	0.8713
7	absolute_error	Log2	10	0.8613
8	absolute_error	Log2	100	0.8722

Best Performing Model:

Random Forest with R² value = 0.8722