

Assignment – Regression Algorithm

1. Problem Statement or Requirement:

Multiple input parameters are provided to predict insurance charges. As the input and output are numeric and the requirement deals with numeric values, selecting the **Machine Learning domain**.

Looking at the dataset, the following are observed:

- Requirement is clear with input(s) and Output
- Independent and Dependent variable are present.

Based on the above observation, **Supervised Learning** is adopted.

The output values are numeric hence proceeding with **Regression**.

2. Basic Info about dataset:

The dataset has multiple inputs such as age, gender, children, smoker, bmi and output charges (insurance charges)...total 1338 rows x 6 columns

3. The inputs such as gender and smoker columns hold the categorical (nominal data) hence these values are converted into numeric using the get_dummies.

4. Using Machine Learning Algorithm, created the following:

- Multiple Linear Regression
- Support Vector Machine
- Decision Tree
- Random Forest

5. Accuracy of the Research Models:

- Multiple Linear Regression: $r^2_score = 1.0$
- Support Vector Machine (without standardization): $r^2_score = 0.9999999999605653$
- Decision Tree (with default parameters): $r^2_score = 0.9982364920332555$
- Random Forest: $r^2_score = 0.999679694289922$

6. The final model selected was Multiple Linear regression with 100% Accuracy when compared to the other models.

SVM(after Standardization)

Hyper Parameter	Linear (r Value)	Rbf(Non Linear)(r Value)	Poly(r Value)	Sigmoid(r Value)
C=1	0.1008799457 2886794	-0.080825646 0270218	-0.041298004 53972243	-0.063665375 37679862
C=10	0.8529720560 922177	-0.008006791 942225089	0.3355578386 035505	0.1488034906 1986875
C=100	0.9999999999 619384	0.4627044357 461626	0.9086496487 205619	0.7958333486 9026
C=500	0.9999999999 802569	0.9072503636 707163	0.9423885132 46426	0.6082403582 150155
C=1000	0.9999999999 814986	0.9635426982 159025	0.9541513072 972456	-0.086276648 21740333
C=2000	0.9999999999 825486	0.9858899155 822013	0.9710280758 19288	-2.667309996 4651152
C=3000	0.9999999999 857418	0.9933146265 127804	0.9822555307 105733	-6.458271017 41801

Decision Tree:

Criterion	Splitter	max_depth	max_leaf_nodes	min_samples_split	min_samples_leaf	R score
Friedman_mse	best	int	int	min	min	0.9989409 919072885
Absolute_error	Best	Int	int	min	min	0.9996490 574735933
poisson	best	int	int	min	min	0.9996835 23863742
Friedman_mse	random	int	int	min	min	0.9994120 464222314
Absolute_error	random	int	int	min	min	0.9980370 731835224
poisson	random	int	int	min	min	0.9985792 862032016

Random Forest

N_estimators	Criterion	Max_features	R_score
50	Absolute error	Sqrt	0.9966520724500993
50	Friedman_mse	sqrt	0.9957778930918343
50	poisson	sqrt	0.996078065629805
50	Absolute error	Log2	0.9966520724500993
50	Friedman_mse	Log2	0.9957778930918343
50	poisson	Log2	0.996078065629805