

Assignment

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.

As dataset scientist, you must develop a model which will predict the insurance charges.

1) Identify your problem statement

Stage 1: Machine learning

Stage 2: Supervised learning

Stage 3: Regression

2) Tell basic info about the dataset (Total number of rows, columns)

Total number of rows are 1338 and the total number of columns is 6

3) Mention the pre-processing method if you're doing any (like converting string to number-nominal data)

It's a nominal dataset, because we can compare the age with others.

4) Develop a good model with r^2 score. You can use any machine learning algorithm; you can create many models. Finally, you have to come up with final model.

My finalized model for this dataset is "Random forest" because the r^2 value of this model is the highest one.

5) All the research value (r^2 score of the models) should be documented.

1) Decision tree

<i>Criterion</i>	<i>Splitter</i>	<i>R2value</i>
Squared_error	<i>best</i>	0.65884579
Friedman_mse	<i>Best</i>	0.6875130
<i>Absolute_error</i>	<i>Best</i>	0.6856941
Max_depth=none	<i>min_samples_split=2</i>	0.69820104
<i>min_samples_leaf=1</i>	<i>min_weight_fraction_leaf=0.0</i>	0.667875
<i>max_features=None</i>	<i>random_state=None</i>	0.69805482
<i>max_leaf_nodes=None</i>	<i>min_impurity_decrease=0.0</i>	0.66857183
<i>ccp_alpha=0.0</i>	<i>monotonic_cst=None</i>	0.66857183
max_depth=3,	<i>min_impurity_decrease=0.4</i>	0.70183154

<i>criterion='squared_error'</i>	<i>random_state=2</i>	0.7015151
<i>criterion='squared_error'</i>	<i>min_samples_split=2</i>	0.68342783
<i>max_depth=None</i>	<i>splitter='best'</i>	0.6884023
<i>min_impurity_decrease=0.0</i>	<i>monotonic_cst=None</i>	0.70506885
<i>ccp_alpha=0.0</i>	<i>min_samples_leaf=1</i>	0.70506885
max_depth=3	min_samples_leaf=100	0.7313303
<i>min_weight_fraction_leaf=0</i>	<i>min_samples_leaf=10</i>	0.87467747
<i>max_leaf_nodes=None</i>	<i>random_state=0</i>	0.690923077

2) SVM

Parameter	parameter	R2 value
	Kernel=Linear	-0.0101954
<i>epsilon=0.0</i>	<i>tol=0.0001</i>	-0.0834051

epsilon=0.0	max_iter=1000	-0.08340516
max_iter=1000	C=1.0	-0.083405160
epsilon=0.0	C=1.0	-0.083405160
epsilon=0.0	verbose=0	-0.08340516
max_iter=1000	C=1.0	-0.0884616
	Kernal='rbf'	-0.0834051
	Kernal='poly'	-0.0757173
	Kernal='sigmoid'	-0.0754463
	gamma='scale'	0.84752647
C=1.0	epsilon=0.2	-0.08340516
C=100	Epsilon=0.6	-0.08340516
C=1000	Epsilon=0.10	0.810719570

C=0.01	Epsilon=100	-0.08868556
	Gamma='auto'	-0.083405160

3) Random forest

<i>n_estimators=50</i>	<i>max_features=1.0</i>	0.8563508
<i>ccp_alpha=0.0</i>	<i>max_features=1.0</i>	0.8475264
<i>min_weight_fraction_leaf=0.0</i>	<i>max_features=1.0</i>	0.85041311
<i>min_weight_fraction_leaf=0.50</i>	<i>max_features=0.20</i>	0.00320632
<i>criterion='squared_error'</i>	<i>min_samples_split=2</i>	0.85108598
<i>criterion='squared_error'</i>	<i>min_samples_split=7</i>	0.872655517
<i>criterion='squared_error'</i>	<i>min_samples_split=10</i>	0.875290100

<i>bootstrap=True</i>	min_samples_split=10	0.87506315
<i>bootstrap=True</i>	min_samples_split=100	0.86859936
<i>max_depth=None</i>	ccp_alpha=0.0	0.85151773
<i>max_depth=None</i>	ccp_alpha=0.20	0.852772384
<i>max_depth=None</i>	monotonic_cst=None	0.84871504
<i>n_jobs=None</i>	n_estimators=100	0.852214440
<i>n_jobs=None</i>	n_estimators=1200	0.853187652

6) Mention your final model and justify why you chosen the same.

Random forest is the best method for the given model , because when we combine these two parameter (*criterion=squared_error,min_samples_split=10*) we get the highest r2 value and the r2 value is 0.875290100.