

In [145]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
# to visualise all the columns in the dataframe
pd.pandas.set_option('display.max_columns', None)
```

In [146]:

```
train_data= pd.read_csv("train_data.csv")
```

In [147]:

```
train_data.head()
```

Out[147]:

	Id	SalePrice	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	LotConfig	LandSlop
0	1	12.247694	0.235294	0.75	0.418208	0.366344	1.0	1.0	0.000000	0.333333	1.0	0.00	0.
1	2	12.109011	0.000000	0.75	0.495064	0.391317	1.0	1.0	0.000000	0.333333	1.0	0.50	0.
2	3	12.317167	0.235294	0.75	0.434909	0.422359	1.0	1.0	0.333333	0.333333	1.0	0.00	0.
3	4	11.849398	0.294118	0.75	0.388581	0.390295	1.0	1.0	0.333333	0.333333	1.0	0.25	0.
4	5	12.429216	0.235294	0.75	0.513123	0.468761	1.0	1.0	0.333333	0.333333	1.0	0.50	0.

In [148]:

```
train_data.shape
```

Out[148]:

(1460, 84)

In [149]:

```
## Capture the dependent feature
y_train=train_data[['SalePrice']]
```

In [150]:

```
## drop dependent feature from dataset
X_train=train_data.drop(['Id','SalePrice'],axis=1)
```

In [151]:

```
### Apply Feature Selection
# first, I specify the Lasso Regression model, and I
# select a suitable alpha (equivalent of penalty).
# The bigger the alpha the less features that will be selected.

# Then I use the selectFromModel object from sklearn, which
# will select the features which coefficients are non-zero

from sklearn.linear_model import Lasso
from sklearn.feature_selection import SelectFromModel
feature_sel_model = SelectFromModel(Lasso(alpha=0.005, random_state=0)) # remember to set the
seed, the random state in this function
feature_sel_model.fit(X_train, y_train)
```

Out[151]:

```
SelectFromModel(estimator=Lasso(alpha=0.005, copy_X=True, fit_intercept=True,
```

```

selection_model(estimator=lasso(alpha=0.0001, copy_X=True, fit_intercept=True,
                                max_iter=1000, normalize=False, positive=False,
                                precompute=False, random_state=0,
                                selection='cyclic', tol=0.0001,
                                warm_start=False),
                max_features=None, norm_order=1, prefit=False, threshold=None)

```

In [152]:

```
feature_sel_model.get_support()
```

Out[152]:

```

array([ True,  True, False, False, False, False, False, False, False,
        False, False,  True, False, False, False, False,  True, False,
        False,  True,  True, False, False, False, False, False, False,
        False, False,  True, False,  True, False, False, False, False,
        False, False, False,  True,  True, False,  True, False, False,
         True,  True, False, False, False, False, False,  True, False,
        False,  True,  True,  True, False,  True,  True, False, False,
        False,  True, False, False, False, False, False, False, False,
        False, False, False, False, False, False,  True, False, False,
        False])

```

In [153]:

```

# let's print the number of total and selected features

# this is how we can make a list of the selected features
selected_feat = X_train.columns[(feature_sel_model.get_support())]

# let's print some stats
print('total features: {}'.format((X_train.shape[1])))
print('selected features: {}'.format(len(selected_feat)))
print('features with coefficients shrank to zero: {}'.format(np.sum(feature_sel_model.estimator_.coef_ == 0)))

```

```

total features: 82
selected features: 21
features with coefficients shrank to zero: 61

```

In [154]:

```
selected_feat
```

Out[154]:

```

Index(['MSSubClass', 'MSZoning', 'Neighborhood', 'OverallQual', 'YearRemodAdd',
       'RoofStyle', 'BsmtQual', 'BsmtExposure', 'HeatingQC', 'CentralAir',
       '1stFlrSF', 'GrLivArea', 'BsmtFullBath', 'KitchenQual', 'Fireplaces',
       'FireplaceQu', 'GarageType', 'GarageFinish', 'GarageCars', 'PavedDrive',
       'SaleCondition'],
      dtype='object')

```

In [155]:

```
X_train=X_train[selected_feat]
```

In [156]:

```
X_train.head()
```

Out[156]:

	MSSubClass	MSZoning	Neighborhood	OverallQual	YearRemodAdd	RoofStyle	BsmtQual	BsmtExposure	HeatingQC	CentralAir
0	0.235294	0.75	0.636364	0.666667	0.883333	0.0	0.75	0.25	1.00	1.0
1	0.000000	0.75	0.500000	0.555556	0.433333	0.0	0.75	1.00	1.00	1.0
2	0.235294	0.75	0.636364	0.666667	0.866667	0.0	0.75	0.50	1.00	1.0

3	0.294118	0.75	0.727273	0.666667	0.333333	0.0	0.50	0.25	0.75	1.0
MSSubClass	MSZoning	Neighborhood	OverallQual	YearRemodAdd	RoofStyle	BsmtQual	BsmtExposure	HeatingQC	CentralAir	
4	0.235294	0.75	1.000000	0.777778	0.833333	0.0	0.75	0.75	1.00	1.0

In [157]:

```
X_train.shape
```

Out[157]:

```
(1460, 21)
```

In [158]:

```
# Import Sci-Kit Learn
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import Normalizer
from sklearn.linear_model import LinearRegression, Lasso
from sklearn.svm import SVR
from sklearn.ensemble import RandomForestRegressor, AdaBoostRegressor, GradientBoostingRegressor,
BaggingRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score, mean_squared_error
from sklearn.model_selection import RandomizedSearchCV, cross_val_score, StratifiedKFold,
learning_curve, KFold

# Ensemble Models
from xgboost import XGBRegressor
```

In [159]:

```
# Use train_test_split from sci-kit learn to segment our data into train and a local testset
X_train, X_test, y_train, y_test = train_test_split(X_train, y_train, test_size=0.2)
```

Define Evaluation Metric

Submissions are evaluated on Root-Mean-Squared-Error (RMSE) between the logarithm of the predicted value and the logarithm of the observed sales price. (Taking logs means that errors in predicting expensive houses and cheap houses will affect the result equally.)

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (P_i - O_i)^2}{n}}$$

We will write a function named rmse to perform this task.

In [160]:

```
def rmse(y_test, y_pred):
    return np.sqrt(mean_squared_error(np.log(y_test), np.log(y_pred)))
```

Modelling We'll start by building standalone models, validating their performance and picking the right ones. Later we will stack all our models into an ensemble for better accuracy.

I just played with a number of models and ended up picking the following models which gave me best results personally. If you find a better way please let me know :)

I tuned the hyperparameters by manually experimenting a lot based on previous experiences, saving you a bunch of time hopefully. Anyways if you find a set of parameters that can work even better let me know in the comments down below.

Linear Regression

Linear Regression

In [161]:

```
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LinearRegression

lin_regressor=LinearRegression()

lin_regressor.fit(X_train, y_train)
mse=cross_val_score(lin_regressor,X_train,y_train,scoring='neg_mean_squared_error',cv=5)
mean_mse=np.mean(mse)
print(mean_mse)
```

-0.018713067659747264

In [162]:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X_train, y_train, test_size=0.2, random_state=0
)
```

In [163]:

```
prediction_linear = lin_regressor.predict(X_test)
```

In [164]:

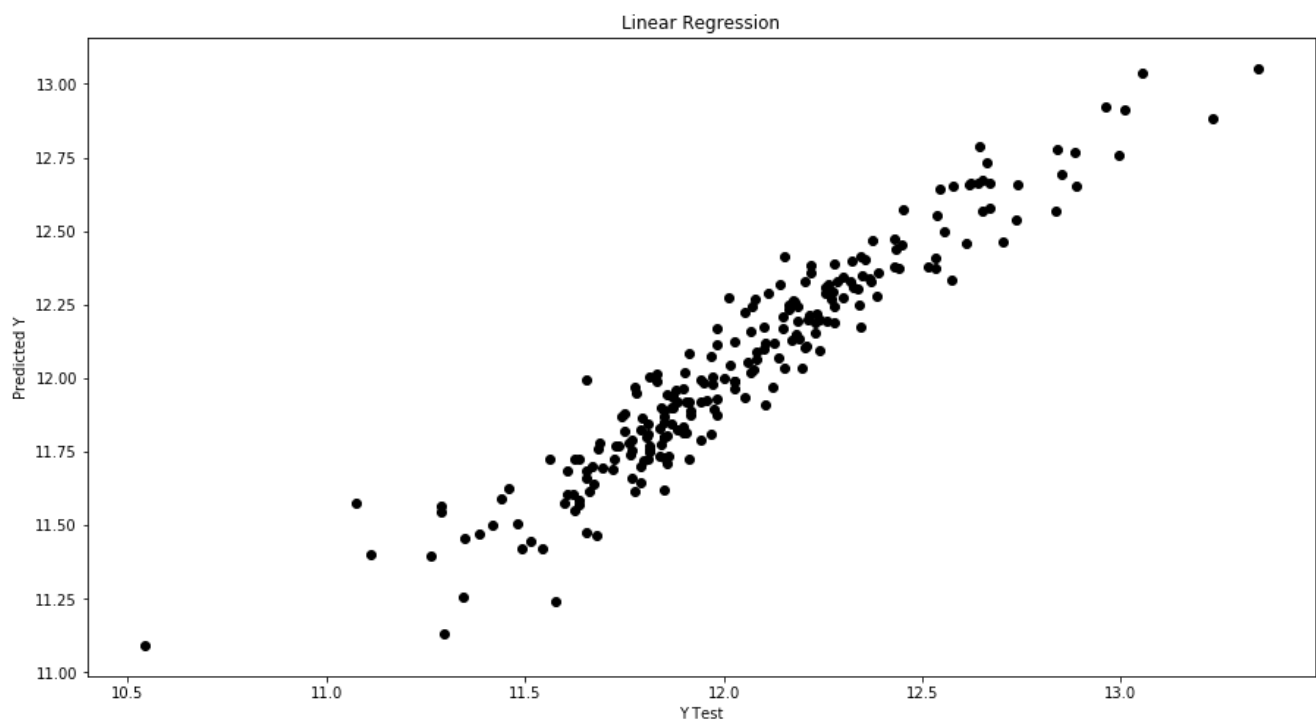
```
rmse(y_test, prediction_linear)
```

Out[164]:

0.010320438372497276

In [165]:

```
plt.figure(figsize=(15,8))
plt.scatter(y_test,prediction_linear, c= 'black')
plt.title("Linear Regression")
plt.xlabel('Y Test')
plt.ylabel('Predicted Y')
plt.show()
```



Ridge Regression

In [166]:

```
from sklearn.linear_model import Ridge
from sklearn.model_selection import GridSearchCV

ridge=Ridge()
parameters={'alpha':[1e-15,1e-10,1e-8,1e-3,1e-2,1,5,10,20,30,35,40,45,50,55,100]}
ridge_regressor=GridSearchCV(ridge,parameters,scoring='neg_mean_squared_error',cv=5)
ridge_regressor.fit(X_train,y_train)
```

Out[166]:

```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=Ridge(alpha=1.0, copy_X=True, fit_intercept=True,
                             max_iter=None, normalize=False, random_state=None,
                             solver='auto', tol=0.001),
             iid='warn', n_jobs=None,
             param_grid={'alpha': [1e-15, 1e-10, 1e-08, 0.001, 0.01, 1, 5, 10,
                                    20, 30, 35, 40, 45, 50, 55, 100]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
             scoring='neg_mean_squared_error', verbose=0)
```

In [167]:

```
print(ridge_regressor.best_params_)
print(ridge_regressor.best_score_)
```

```
{'alpha': 0.01}
-0.01973394565423983
```

In [168]:

```
prediction_ridge=ridge_regressor.predict(X_test)
```

In [169]:

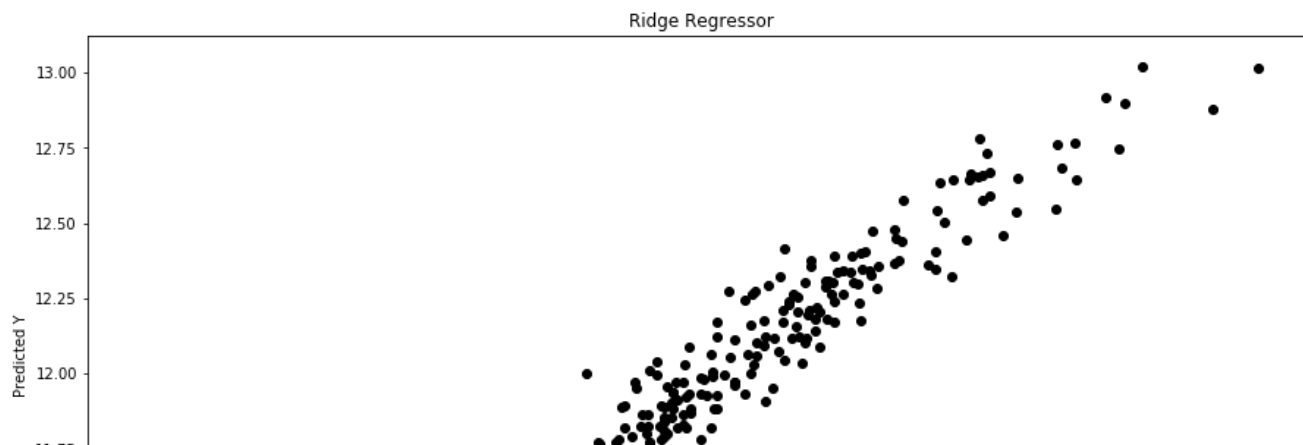
```
rmse(y_test, prediction_ridge)
```

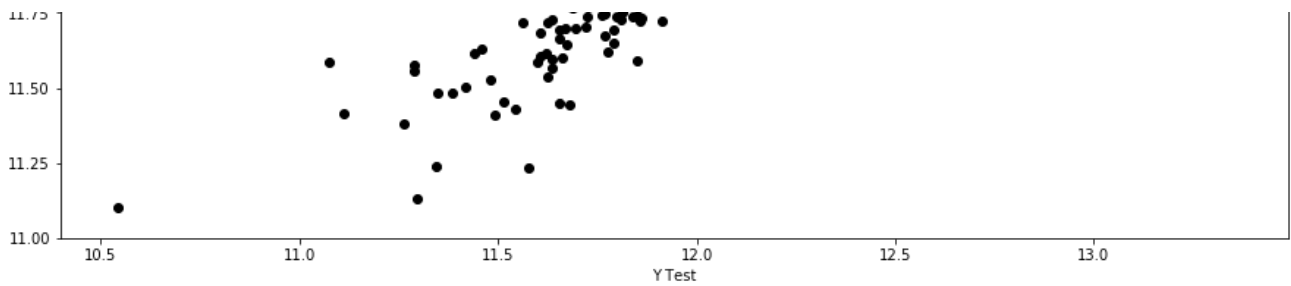
Out[169]:

```
0.010653454462053609
```

In [170]:

```
plt.figure(figsize=(15,8))
plt.scatter(y_test,prediction_ridge, c= 'black')
plt.title("Ridge Regressor")
plt.xlabel('Y Test')
plt.ylabel('Predicted Y')
plt.show()
```





Lasso Regression

In [171]:

```
from sklearn.linear_model import Lasso
from sklearn.model_selection import GridSearchCV
lasso=Lasso()
parameters={'alpha':[1e-15,1e-10,1e-8,1e-3,1e-2,1,5,10,20,30,35,40,45,50,55,100]}
lasso_regressor=GridSearchCV(lasso,parameters,scoring='neg_mean_squared_error',cv=5)

lasso_regressor.fit(X_train,y_train)
print(lasso_regressor.best_params_)
print(lasso_regressor.best_score_)
```

```
{'alpha': 1e-08}
-0.019735621986625127
```

In [172]:

```
prediction_lasso=lasso_regressor.predict(X_test)
```

In [173]:

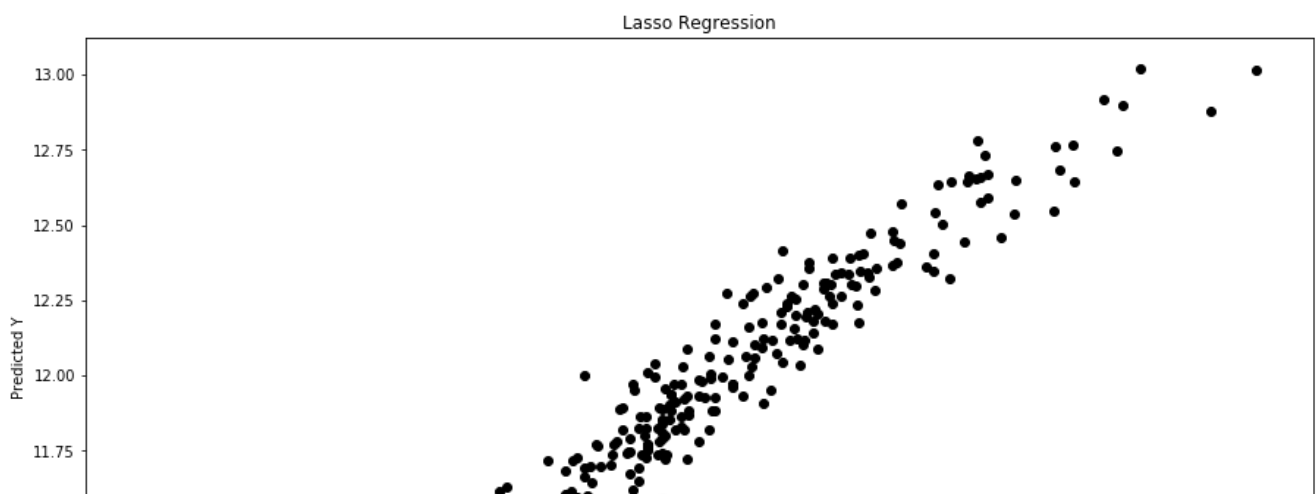
```
rmse(y_test, prediction_lasso)
```

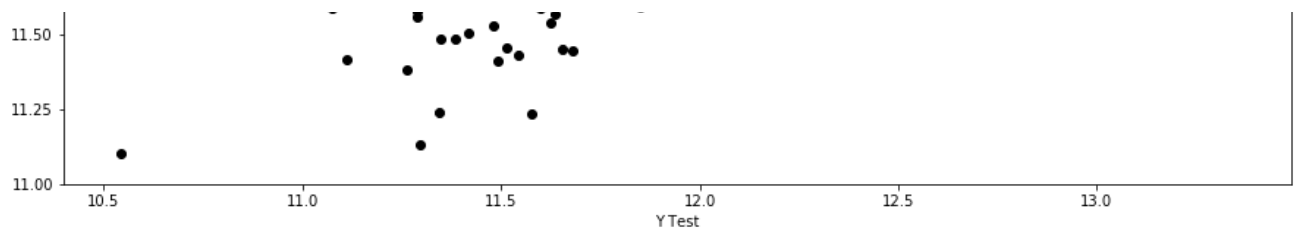
Out[173]:

```
0.010651256824970728
```

In [174]:

```
plt.figure(figsize=(15,8))
plt.scatter(y_test,prediction_lasso, c= 'black')
plt.title("Lasso Regression")
plt.xlabel('Y Test')
plt.ylabel('Predicted Y')
plt.show()
```





In [175]:

```
#Train the model
from sklearn import linear_model
model = linear_model.LinearRegression()
```

In [176]:

```
#Fit the model
model.fit(X_train, y_train)
```

Out[176]:

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

In [177]:

```
linear_pred = model.predict(X_test)
```

In [178]:

```
score = model.score(X_test, y_test)
print(score)
```

```
0.9004020643396348
```

In [179]:

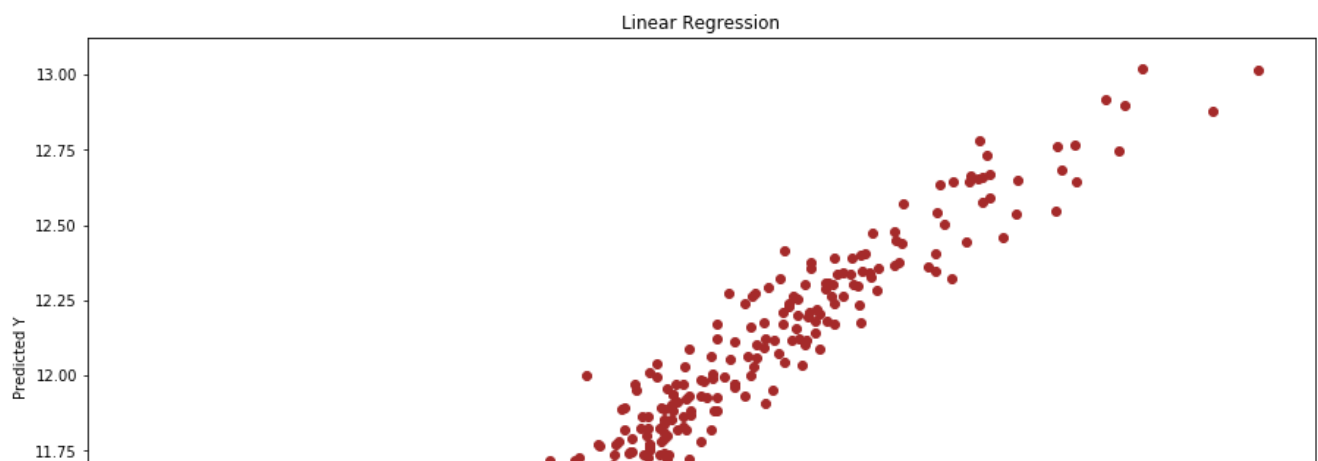
```
rmse(y_test, linear_pred)
```

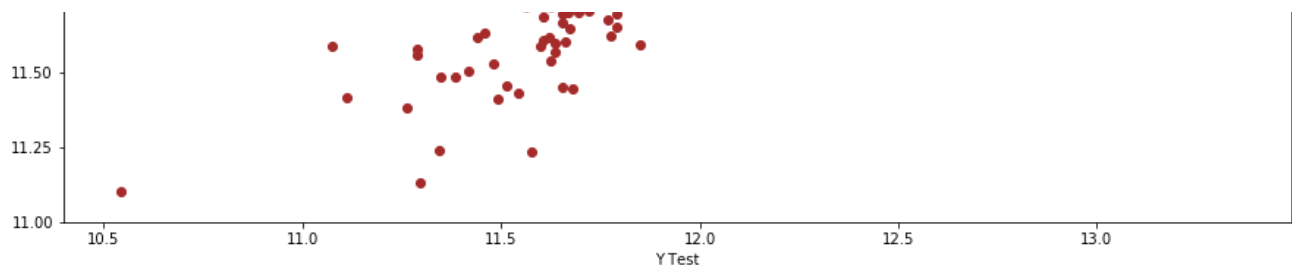
Out[179]:

```
0.01065125433107904
```

In [183]:

```
plt.figure(figsize=(15,8))
plt.scatter(y_test, linear_pred, c= 'brown')
plt.title("Linear Regression")
plt.xlabel('Y Test')
plt.ylabel('Predicted Y')
plt.show()
```





Random Forest

A Random Forest is an ensemble technique capable of performing both regression and classification tasks with the use of multiple decision trees and a technique called Bootstrap Aggregation, commonly known as bagging. What is bagging you may ask? Bagging, in the Random Forest method, involves training each decision tree on a different data sample where sampling is done with replacement.

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (P_i - O_i)^2}{n}}$$

In [184]:

```
# Initialize the model
random_forest = RandomForestRegressor(n_estimators=1200,
                                     max_depth=15,
                                     min_samples_split=5,
                                     min_samples_leaf=5,
                                     max_features=None,
                                     random_state=42,
                                     oob_score=True
                                     )

# Perform cross-validation to see how well our model does
kf = KFold(n_splits=5)
y_pred = cross_val_score(random_forest, X_train, y_train, cv=kf, n_jobs=-1)
y_pred.mean()
```

Out[184]:

0.8509875471915518

In [185]:

```
# Fit the model to our data
random_forest.fit(X_train, y_train)
```

C:\Users\Mahesh\Anaconda3\lib\site-packages\ipykernel_launcher.py:2: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

Out[185]:

```
RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=15,
                      max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=5, min_samples_split=5,
                      min_weight_fraction_leaf=0.0, n_estimators=1200,
                      n_jobs=None, oob_score=True, random_state=42, verbose=0,
                      warm_start=False)
```

In [186]:

```
# Make predictions on test data
```


[illegible]

```
        scale_pos_weight=1, seed=27,  
        reg_alpha=0.00006  
    )
```

```
# Perform cross-validation to see how well our model does  
kf = KFold(n_splits=5)  
y_pred = cross_val_score(xg_boost, X_train, y_train, cv=kf, n_jobs=-1)  
y_pred.mean()
```

Out[190]:

0.8459129219428562

In [191]:

```
# Fit our model to the training data  
xg_boost.fit(X_train,y_train)
```

[09:47:26] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

Out[191]:

```
XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,  
             colsample_bynode=1, colsample_bytree=0.2, gamma=0.6,  
             importance_type='gain', learning_rate=0.01, max_delta_step=0,  
             max_depth=4, min_child_weight=1, missing=None, n_estimators=6000,  
             n_jobs=1, nthread=-1, objective='reg:linear', random_state=0,  
             reg_alpha=6e-05, reg_lambda=1, scale_pos_weight=1, seed=27,  
             silent=None, subsample=0.7, verbosity=1)
```

In [192]:

```
# Make predictions on the test data  
xgb_pred = xg_boost.predict(X_test)
```

In [193]:

```
rmse(y_test, xgb_pred)
```

Out[193]:

0.011868828262320705

In [194]:

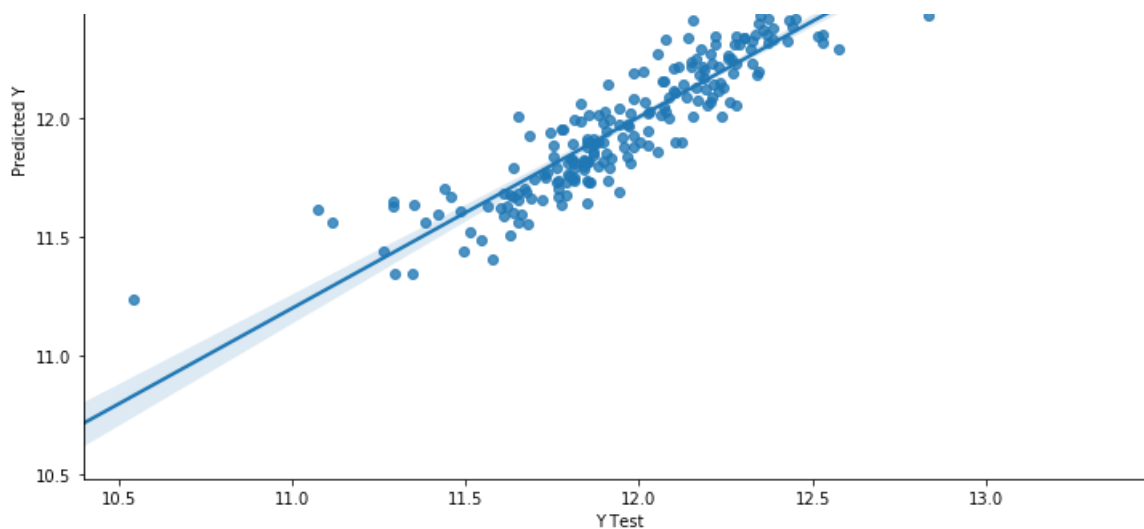
```
score = xg_boost.score(X_test, y_test)  
print(score)
```

0.8763628666070651

In [195]:

```
plt.figure(figsize=(12,8))  
  
sns.regplot(x=y_test,y=xgb_pred,truncate=False)  
plt.xlabel('Y Test')  
plt.ylabel('Predicted Y')  
plt.show()
```





Gradient Boost Regressor

Gradient Boosting trains many models in a gradual, additive and sequential manner. The major difference between AdaBoost and Gradient Boosting Algorithm is how the two algorithms identify the shortcomings of weak learners (eg. decision trees). While the AdaBoost model identifies the shortcomings by using high weight data points, gradient boosting performs the same by using gradients in the loss function ($y=ax+b+e$, e needs a special mention as it is the error term). The loss function is a measure indicating how good are model's coefficients are at fitting the underlying data. A logical understanding of loss function would depend on what we are trying to optimise. We are trying to predict the sales prices by using a regression, then the loss function would be based off the error between true and predicted house prices.

In [196]:

```
# Initialize our model
g_boost = GradientBoostingRegressor( n_estimators=6000, learning_rate=0.01,
                                     max_depth=5, max_features='sqrt',
                                     min_samples_leaf=15, min_samples_split=10,
                                     loss='ls', random_state =42
                                   )

# Perform cross-validation to see how well our model does
kf = KFold(n_splits=5)
y_pred = cross_val_score(g_boost, X_train, y_train, cv=kf, n_jobs=-1)
y_pred.mean()
```

Out[196]:

0.8801162764904106

In [197]:

```
# Fit our model to the training data
g_boost.fit(X_train,y_train)
```

C:\Users\Mahesh\Anaconda3\lib\site-packages\sklearn\ensemble\gradient_boosting.py:1450:
DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
y = column_or_1d(y, warn=True)

Out[197]:

```
GradientBoostingRegressor(alpha=0.9, criterion='friedman_mse', init=None,
                           learning_rate=0.01, loss='ls', max_depth=5,
                           max_features='sqrt', max_leaf_nodes=None,
                           min_impurity_decrease=0.0, min_impurity_split=None,
                           min_samples_leaf=15, min_samples_split=10,
                           min_weight_fraction_leaf=0.0, n_estimators=6000,
                           n_iter_no_change=None, presort='auto',
                           random_state=42, subsample=1.0, tol=0.0001,
                           validation_fraction=0.1, verbose=0, warm_start=False)
```

In [198]:

```
# Make predictions on test data
gbm_pred = g_boost.predict(X_test)
```

In [199]:

```
rmse(y_test, gbm_pred)
```

Out[199]:

0.01031406660638029

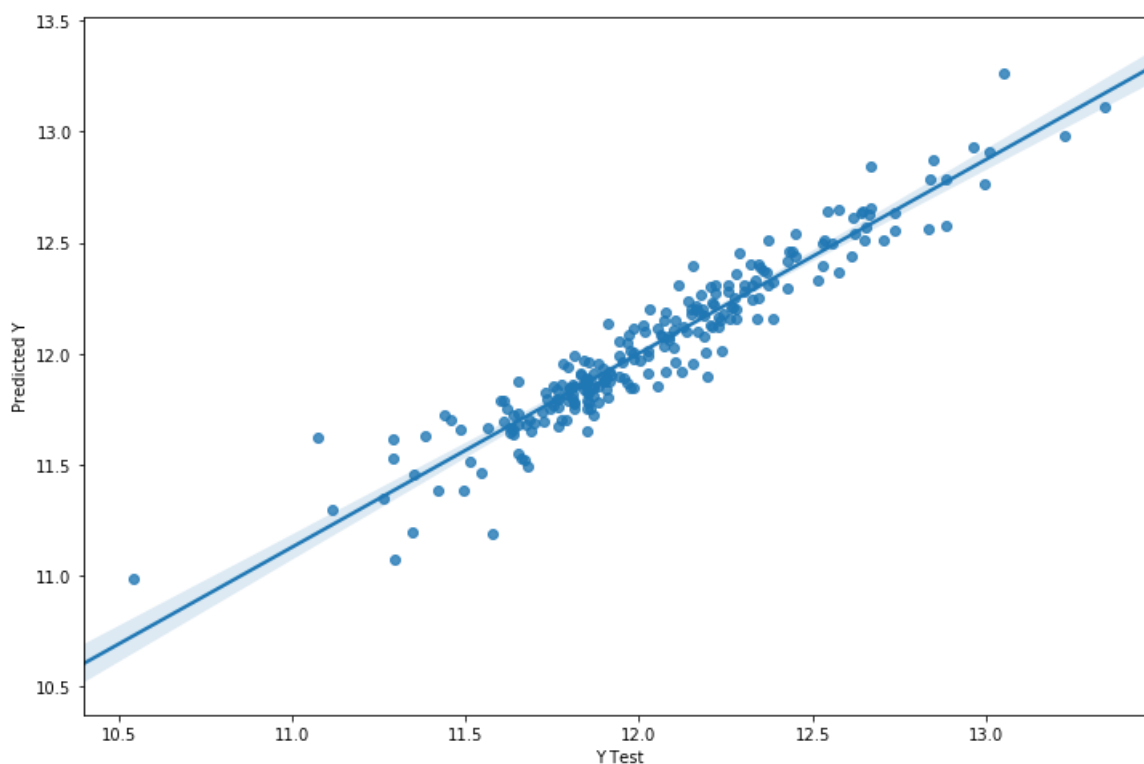
In [200]:

```
score = g_boost.score(X_test, y_test)
print(score)
```

0.906724498774255

In [201]:

```
plt.figure(figsize=(12,8))
sns.regplot(x=y_test,y=gbm_pred,truncate=False)
plt.xlabel('Y Test')
plt.ylabel('Predicted Y')
plt.show()
```



XGBoost With Hyper parameter tuning

In [202]:

```
# List of the models to be stacked
models = [g_boost, xg_boost, random_forest]
```

In [203]:

```
import xgboost
```

```
classifier=xgboost.XGBRegressor()
```

In [204]:

```
import xgboost
regressor=xgboost.XGBRegressor()
```

In [205]:

```
booster=['gbtree','gblinear']
base_score=[0.25,0.5,0.75,1]
```

In [206]:

```
## Hyper Parameter Optimization

n_estimators = [100, 500, 900, 1100, 1500]
max_depth = [2, 3, 5, 10, 15]
booster=['gbtree','gblinear']
learning_rate=[0.05,0.1,0.15,0.20]
min_child_weight=[1,2,3,4]

# Define the grid of hyperparameters to search
hyperparameter_grid = {
    'n_estimators': n_estimators,
    'max_depth':max_depth,
    'learning_rate':learning_rate,
    'min_child_weight':min_child_weight,
    'booster':booster,
    'base_score':base_score
}
```

In [207]:

```
# Set up the random search with 4-fold cross validation
random_cv = RandomizedSearchCV(estimator=regressor,
    param_distributions=hyperparameter_grid,
    cv=5, n_iter=50,
    scoring = 'neg_mean_absolute_error',n_jobs = 4,
    verbose = 5,
    return_train_score = True,
    random_state=42)
```

In [208]:

```
random_cv.fit(X_train,y_train)
```

Fitting 5 folds for each of 50 candidates, totalling 250 fits

```
[Parallel(n_jobs=4)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=4)]: Done 10 tasks      | elapsed:    2.1s
[Parallel(n_jobs=4)]: Done 64 tasks      | elapsed:   21.4s
[Parallel(n_jobs=4)]: Done 154 tasks     | elapsed:   38.1s
[Parallel(n_jobs=4)]: Done 250 out of 250 | elapsed:   55.0s finished
```

[09:48:38] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

Out[208]:

```
RandomizedSearchCV(cv=5, error_score='raise-deprecating',
    estimator=XGBRegressor(base_score=0.5, booster='gbtree',
        colsample_bylevel=1,
        colsample_bynode=1,
        colsample_bytree=1, gamma=0,
        importance_type='gain',
        learning_rate=0.1, max_delta_step=0,
        max_depth=3, min_child_weight=1,
        missing=None, n_estimators=100,
```

```

        n_jobs=1, nthread=None,
        objective='reg:linear',
        random_st...
iid='warn', n_iter=50, n_jobs=4,
param_distributions={'base_score': [0.25, 0.5, 0.75, 1],
                    'booster': ['gbtree', 'gblinear'],
                    'learning_rate': [0.05, 0.1, 0.15, 0.2],
                    'max_depth': [2, 3, 5, 10, 15],
                    'min_child_weight': [1, 2, 3, 4],
                    'n_estimators': [100, 500, 900, 1100,
                                     1500]},
pre_dispatch='2*n_jobs', random_state=42, refit=True,
return_train_score=True, scoring='neg_mean_absolute_error',
verbose=5)

```

In [209]:

```
random_cv.best_estimator_
```

Out [209]:

```

XGBRegressor(base_score=1, booster='gbtree', colsample_bylevel=1,
             colsample_bynode=1, colsample_bytree=1, gamma=0,
             importance_type='gain', learning_rate=0.15, max_delta_step=0,
             max_depth=5, min_child_weight=2, missing=None, n_estimators=100,
             n_jobs=1, nthread=None, objective='reg:linear', random_state=0,
             reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
             silent=None, subsample=1, verbosity=1)

```

In [210]:

```

regressor=xgboost.XGBRegressor(base_score=0.25, booster='gbtree', colsample_bylevel=1,
                               colsample_bynode=1, colsample_bytree=1, gamma=0,
                               importance_type='gain', learning_rate=0.05, max_delta_step=0,
                               max_depth=2, min_child_weight=4, missing=None, n_estimators=900,
                               n_jobs=1, nthread=None, objective='reg:linear', random_state=0,
                               reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
                               silent=None, subsample=1, verbosity=1)

```

In [211]:

```
regressor.fit(X_train,y_train)
```

[09:48:39] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

Out [211]:

```

XGBRegressor(base_score=0.25, booster='gbtree', colsample_bylevel=1,
             colsample_bynode=1, colsample_bytree=1, gamma=0,
             importance_type='gain', learning_rate=0.05, max_delta_step=0,
             max_depth=2, min_child_weight=4, missing=None, n_estimators=900,
             n_jobs=1, nthread=None, objective='reg:linear', random_state=0,
             reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
             silent=None, subsample=1, verbosity=1)

```

In [212]:

```
xgb_pred = regressor.predict(X_test)
```

In [213]:

```
rmse(y_test, xgb_pred)
```

Out [213]:

```
0.009900999274032164
```

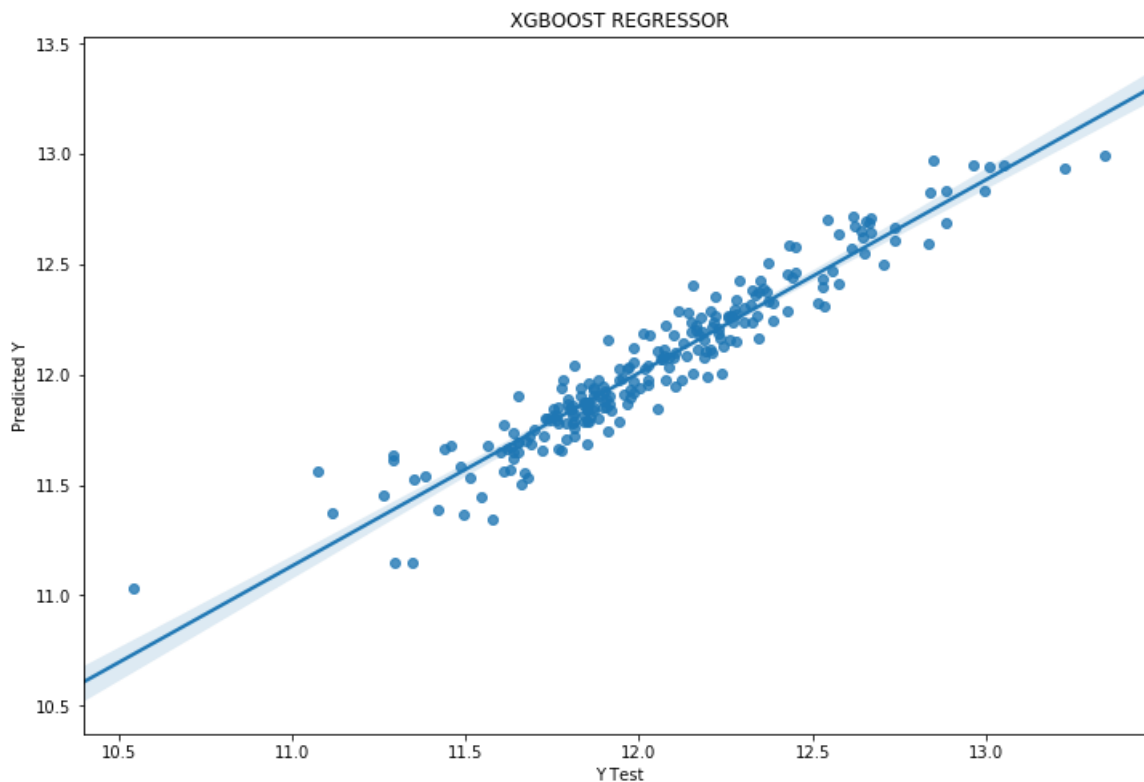
In [214]:

```
score = regressor.score(X_test, y_test)
print(score)
```

0.9141605278376341

In [215]:

```
plt.figure(figsize=(12,8))
sns.regplot(x=y_test,y=xgb_pred,truncate=False)
plt.title("XGBOOST REGRESSOR")
plt.xlabel('Y Test')
plt.ylabel('Predicted Y')
plt.show()
```



Decision Tree Regressor

The decision tree is a simple machine learning model for getting started with regression tasks.

Background A decision tree is a flow-chart-like structure, where each internal (non-leaf) node denotes a test on an attribute, each branch represents the outcome of a test, and each leaf (or terminal) node holds a class label. The topmost node in a tree is the root node. (see here for more details).

In [216]:

```
from sklearn.tree import DecisionTreeRegressor
from sklearn import metrics
dtreg = DecisionTreeRegressor(random_state = 100)
dtreg.fit(X_train, y_train)
```

Out[216]:

```
DecisionTreeRegressor(criterion='mse', max_depth=None, max_features=None,
                      max_leaf_nodes=None, min_impurity_decrease=0.0,
                      min_impurity_split=None, min_samples_leaf=1,
                      min_samples_split=2, min_weight_fraction_leaf=0.0,
                      presort=False, random_state=100, splitter='best')
```

In [220]:

```
dtr_pred = dtreg.predict(X_test)
```

```
In [221]:
```

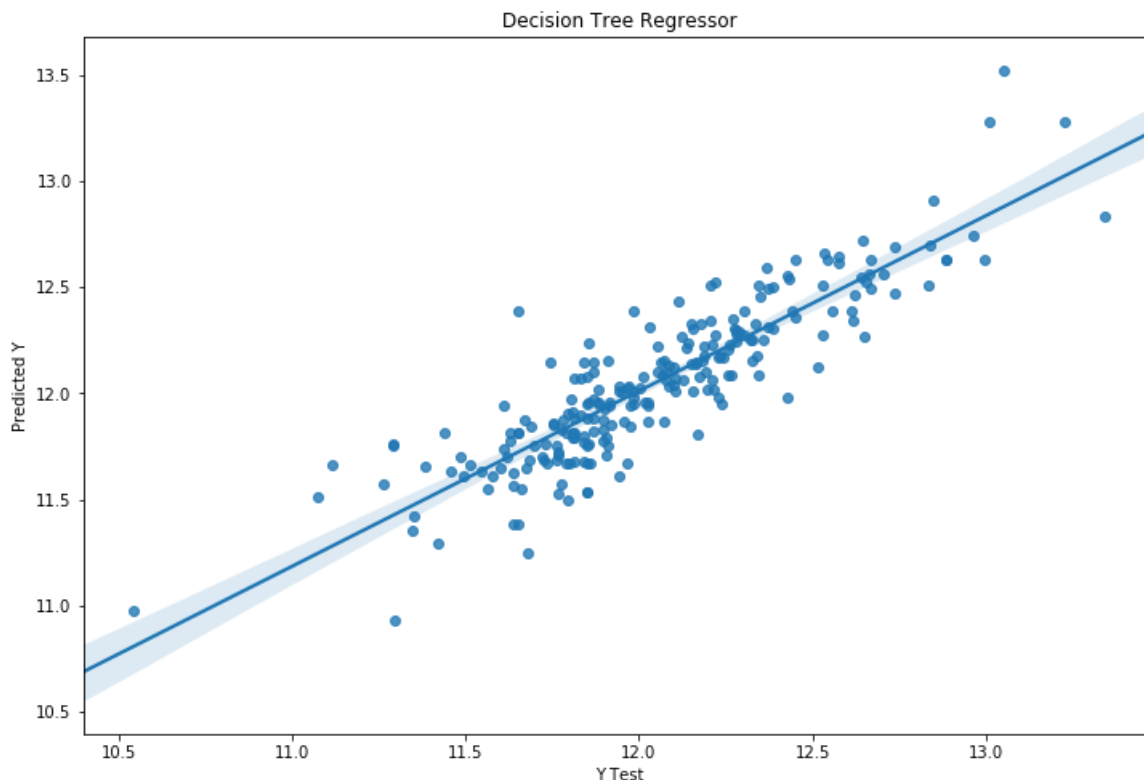
```
rmse(y_test, dtr_pred)
```

```
Out[221]:
```

```
0.01576182946159133
```

```
In [222]:
```

```
plt.figure(figsize=(12,8))  
sns.regplot(x=y_test,y=dtr_pred,truncate=False)  
plt.title("Decision Tree Regressor")  
plt.xlabel('Y Test')  
plt.ylabel('Predicted Y')  
plt.show()
```



Model Comparison

We can say the best working model by looking RMSE rates The best working model is XGBoost Regressor. We are going to see the error rate. which one is better?

```
In [223]:
```

```
error_rate=np.array([rmse(y_test, linear_pred),rmse(y_test, rf_pred),rmse(y_test, dtr_pred),rmse(y_test, gbm_pred),rmse(y_test, xgb_pred)])
```

```
In [228]:
```

```
error_rate
```

```
Out[228]:
```

```
array([0.01065125, 0.01129651, 0.01576183, 0.01031407, 0.009901  ])
```

```
In [232]:
```



```

from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Model", "Error"]

x.add_row(["XG Boost", 0.009901])
x.add_row(["Gradient Boost", 0.01031407])
x.add_row(["Linear Regression", 0.01065125])
x.add_row(["Random Forest", 0.01129651 ])
x.add_row(["Decision Tree", 0.01576183])
print(x)

```

```

+-----+-----+
|      Model      |   Error   |
+-----+-----+
|      XG Boost   | 0.009901  |
| Gradient Boost  | 0.01031407|
| Linear Regression| 0.01065125|
| Random Forest   | 0.01129651|
| Decision Tree   | 0.01576183|
+-----+-----+

```

In the above table, we compare the different models and its RMSE values. And we observed that XG Boost gives the best value when compared to all other models.

In []: