

Background on Data

The dataset relating to this competition considers features pertaining buying a home. The data being used is a compilation of 79 features, all of which describe almost every aspect of homes in Ames, Iowa. The goal of this competition is to use these given features as a means of prediction to project the final price of each home.

Data Preparation

From a data preparation standpoint, I first retrieved the column names, as the dataset contained a vast number of features. I then wanted to see which features had null values by the creating a bar plot that represented the features that had the greatest number of absent values in both the training and test sets. From this point, I felt it to be necessary to create a list, containing only features with missing values and used keys ('MissingValues' and 'Fraction') to display the number of null instances and the percentage of missing data per feature. Following this step, a correlation matrix was then generated to determine which features had the biggest influence on sale price. This correlation matrix, as it relates to sale price, was then used to filter out features with a correlation greater than or equal the absolute value of 0.6 and then stored as the object 'interesting variables'. I then created my training and test sets using the aforementioned data object.

Results/Evaluation

After creating the training and test datasets, I fit the data to ridge, elastic net and lasso regressors in which I used a cross validation design that incorporated the root mean-squared error (RMSE). From the regressors listed, lasso regression performed the best with a RMSE of 44,998.85

followed by ridge regression (RMSE = 56,373.32). Subsequent to this cross-validation design, I then decided to employ random forest methods, one of which used gradient boost. The RMSE for the random forest was 46,000.41 while the score for the gradient boosted random forest came to be 38,915.38.

Model Performance and Management Recommendation

As it pertains to model performance, Kaggle favored the gradient boosted model as it resulted in a score of 0.20640 while the score for the random forest model without boost came to be 0.26096. From a managerial perspective, the gradient boosted random forest model should be used when assessing the market value of residential real estate as proven from the RMSE scores of both respective models. The score is indicative of the performance of the model on data that it has yet to see, hence the model with the best score should be used. The model with the best score in this case would be the gradient boosted random forest.

Appendix


```
In [ ]: from google.colab import drive
drive.mount('/content/drive')
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount('/c
ontent/drive', force_remount=True).
```

```
In [ ]: import os
import pandas as pd
import seaborn as sns
import numpy as np
import plotly.graph_objs as go
from plotly.offline import download_plotlyjs, init_notebook_mode, plot, iplot
init_notebook_mode(connected=True)
from scipy.stats import norm, skew
from sklearn import tree
from sklearn.metrics import accuracy_score
from sklearn import preprocessing
from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import Imputer
from sklearn.preprocessing import MaxAbsScaler
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import RobustScaler
from sklearn.preprocessing import Normalizer
from sklearn.preprocessing import QuantileTransformer
from sklearn.preprocessing import PowerTransformer
matplotlib inline
import matplotlib as mpl
from matplotlib import cm
import matplotlib.pyplot as plt
from IPython.core.interactiveshell import InteractiveShell
from sklearn.linear_model import RidgeCV, ElasticNetCV, LassoCV
from sklearn.decomposition import PCA #Principal Component Analysis
from sklearn.preprocessing import LabelEncoder, transform_categorical into numbers
from sklearn.model_selection import KFold
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
from sklearn.ensemble import VotingClassifier, GradientBoostingRegressor
from sklearn.linear_model import LinearRegression, Lasso, Ridge, ElasticNet
from sklearn.svm import SVC
from sklearn.metrics import mean_squared_error, make_scorer
import xgboost
from sklearn.model_selection import GridSearchCV
```

```
In [ ]: housing_train= pd.read_csv('train.csv')
housing_test= pd.read_csv('test.csv')
```

```
In [ ]: housing_train.head()
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Ally	LotShape	LandContour	Utilities	LotConfig	LandSlope	Y
0	1	60	RL	85.0	8450	Pave	NaN	Reg	Lvl	AI Pub	Inside	GE	
1	2	11	RM	60.0	9600	Pave	NaN	IR1	Lvl	AI Pub	FR2	GE	
2	3	4	RM	60.0	11250	Pave	NaN	IR1	Lvl	AI Pub	Inside	GE	
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AI Pub	Cornr	GE	
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AI Pub	FR2	GE	

```
5 rows x 13 columns
```

```
In [ ]: housing_test.head()
```

```
Out [ ]: missing SalePrice
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Ally	LotShape	LandContour	Utilities	LotConfig	LandSlope	Y
0	1461	20	RH	80.0	11822	Pave	NaN	Reg	Lvl	AI Pub	Inside	GE	
1	1462	20	RL	81.0	14297	Pave	NaN	IR1	Lvl	AI Pub	Cornr	GE	
2	1463	60	RL	74.0	13830	Pave	NaN	IR1	Lvl	AI Pub	Inside	GE	
3	1464	60	RL	78.0	9978	Pave	NaN	IR1	Lvl	AI Pub	Inside	GE	
4	1465	120	RL	43.0	5005	Pave	NaN	IR1	HLS	AI Pub	Inside	GE	

```
In [ ]: housing_train.columns
```

```
Out [ ]: Index(['Id', 'MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street', 'Ally', 'LotShape', 'LandContour', 'Utilities', 'LotConfig', 'LandSlope', 'Y'],
      dtype='object', name='Index')
MSSubClass    int64
MSZoning       object
LotFrontage    float64
LotArea        int64
...           ...
MiscVal       int64
YrSold        object
SaleType      object
SaleCondition  object
SalePrice     int64
Length: 81, dtype: object
```

```
In [ ]: print ("Train set: ", housing_train.shape)
print ("Test set: ", housing_test.shape)
```

```
Train set: (1468, 81)
Test set: (1459, 80)
```

```
In [ ]: housing_train.isnull().sum()
```

```
Out [ ]: Id      0
MSSubClass  0
MSZoning     0
LotFrontage 227
LotArea      0
...
MiscVal     0
YrSold      0
SaleType    0
SaleCondition 0
SalePrice   0
Length: 81, dtype: int64
```

```
In [ ]: housing_test.isnull().sum()
```

```
Out [ ]: Id      0
MSSubClass  0
MSZoning     4
LotFrontage 227
LotArea      0
...
MiscVal     0
YrSold      0
SaleType    1
SaleCondition 0
Length: 80, dtype: int64
```

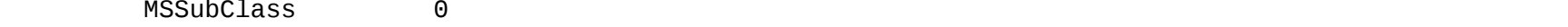
```
In [ ]: housing_test.loc[housing_test.SaleType.isnull()]
```

```
Out [ ]: Id      1469
MSSubClass  20
MSZoning     RL
LotFrontage 20
LotArea      0
...
MiscVal     0
YrSold      0
SaleType    0
SaleCondition 0
Length: 80, dtype: int64
```

```
In [ ]: #Top 20 Features with missing data from housing_train
```

```
sns.set_style('whitegrid')
plt.style.use('fivethirtyeight')
plt.figure(figsize=(15,4))
df=pd.Series(1 - housing_test.count() / len(housing_train)).sort_values(ascending=False).head(20)
sns.barplot(x=df.index, y=df,palette='Blues_d')
plt.xticks(rotation=90)
```

```
Out [ ]: (array([ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19]), <a list of 20 Text major ticklabel objects>)
```



```
In [ ]: #Top 20 Features with missing data from housing_test
```

```
sns.set_style('whitegrid')
plt.style.use('fivethirtyeight')
plt.figure(figsize=(15,4))
df=pd.Series(1 - housing_test.count() / len(housing_test)).sort_values(ascending=False).head(20)
sns.barplot(x=df.index, y=df,palette='Blues_d')
plt.xticks(rotation=90)
```

```
Out [ ]: (array([ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19]), <a list of 20 Text major ticklabel objects>)
```



```
In [ ]: #Houses sold by year housing_train
```

```
housing_train['YrSold'].value_counts()
```

```
Out [ ]: 2009    338
2008    329
2006    314
2009    304
2010    175
Name: YrSold, dtype: int64
```

```
In [ ]: #Houses sold by year housing_test
```

```
housing_test['YrSold'].value_counts()
```

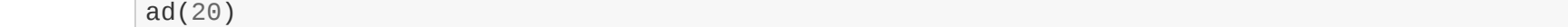
```
Out [ ]: 2007    343
2008    318
2009    289
2006    395
2010    164
Name: YrSold, dtype: int64
```

```
In [ ]: #Housing price = mean sale price by year
```

```
grp_year=housing_train.groupby('YrSold')
plt.figure(figsize=(5,3))
```

```
housing_years=grp_year['SalePrice'].mean().reset_index()
sns.barplot(x=housing_years['YrSold'], y=housing_years['SalePrice'],palette='Blues_d')
plt.xticks(rotation=8)
```

```
Out [ ]: (array([0, 1, 2, 3, 4]), <a list of 5 Text major ticklabel objects>)
```

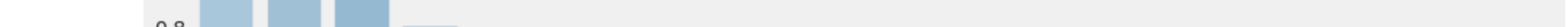


```
In [ ]: #Sales Price Distribution
```

```
plt.figure(figsize=(12,4))
```

```
sns.distplot(housing_train['SalePrice'], fit=norm);
(mu, sigma) = norm.fit(housing_train['SalePrice'])
print( 'mu = %1.2f and sigma = %1.2f'%(mu,sigma))
plt.legend(['Normal Distribution (%mu=%1.2f and %sigma=%1.2f)'.format(mu, sigma),
            'loc = right'])
```

```
Out [ ]: (array([ 0, 1, 2, 3, 4]), <a list of 5 Text major ticklabel objects>)
```



```
In [ ]: housing_train.nona=housing_train.dropna()
print(housing_train.nona.shape)
print(housing_train.shape)
```

```
Out [ ]: (0, 81)
(1468, 81)
```

```
In [ ]: num_missing = housing_train.isnull().sum()
percent = num_missing / housing_train.isnull().count()
```

```
df_missing = pd.concat([num_missing, percent], axis=1, keys=['MissingValues', 'Fraction'])
df_missing=df_missing.sort_values('Fraction', ascending=False)
```

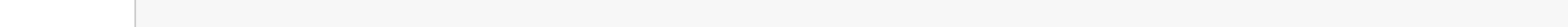
```
Out [ ]: MissingValues Fraction
PoolQC      1463  0.996205
MiscFeature 1406  0.963014
Alley       1369  0.937671
Fence       1179  0.807534
FireplaceQu 690  0.472603
LotFrontage 259  0.176901
GarageYRBL  81  0.055479
GarageCond  81  0.055479
GarageType  81  0.055479
GarageFinH  81  0.055479
GarageQual  81  0.055479
BasementFin2 38  0.026027
BsmtExposure 38  0.026027
BsmtQual    37  0.023432
BsmtCond    37  0.023432
BsmtFinType1 37  0.023432
MaxYrBuilt   8  0.005479
MasVnrType   8  0.005479
Electrical   1  0.000695
```

```
In [ ]: variables_to_keep = df_missing[df_missing['MissingValues'] == 0].index
df_train = housing_train[variables_to_keep]
df_test = housing_test[variables_to_keep]
df_train.shape
```

```
Out [ ]: (1468, 82)
```

```
In [ ]: matrix = df_train.corr()
f, ax = plt.subplots(figsize=(16, 12))
sns.heatmap(matrix, vmax=0.7, square=True)
```

```
Out [ ]: matplotlib.axes._subplots.AxesSubplot at 0x7F80C9a5E1D0
```



```
In [ ]: # Filter out the target variables (SalePrice) and variables with a low correlation score (v
such that <math>\rho < 0.6</math>)
```

```
interesting_variables = matrix['SalePrice'].sort_values(ascending=False)
interesting_variables = interesting_variables[abs(interesting_variables) >= 0.6]
cols = interesting_variables.index.values.tolist()
sns.pairplot(df_train[cols], size=2.0)
```

```
Out [ ]: /usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:1960: UserWarning:
The 'size' parameter has been renamed to 'height'; please update your code.
```



```
In [ ]: pred_vars = {v for v in interesting_variables.index values if v != 'SalePrice'}
target_var = 'SalePrice'
```

```
X = df_train[pred_vars]
y = df_train[target_var]
X_train=X
X_test=X_test
y_train=y_train
y_test=y_test
```

```
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)
```

```
Out [ ]: /usr/local/lib/python3.7/dist-packages/pandas/core/series.py:4536: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide
id=Indexing.html#returning-a-view-versus-a-copy
```

```
Out [ ]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
criterion='gini', max_depth=None, max_features='auto',
max_leaf_nodes=None, max_samples=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=100,
n_jobs=None, oob_score=False, random_state=42, verbose=0,
warm_start=False)
```

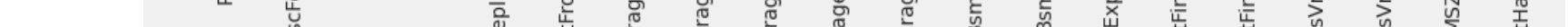
```
In [ ]: num_missing = X_train.isnull().sum()
percent = num_missing / X_train.isnull().count()
```

```
df_missing = pd.concat([num_missing, percent], axis=1, keys=['MissingValues', 'Fraction'])
df_missing=df_missing.sort_values('Fraction', ascending=False)
```

```
Out [ ]: MissingValues Fraction
```

```
In [ ]: y_pred = model.predict(X_test)
plt.scatter(X_test, y_test)
plt.xlabel('Prediction')
plt.ylabel('Real value')
```

```
Out [ ]: # Now add the perfect prediction line
diagonal = np.linspace(0, np.max(X_test), 100)
plt.plot(diagonal, diagonal, '-r')
plt.show()
```



```
In [ ]: np.sqrt(mean_squared_error(y_test, y_pred))
```

```
Out [ ]: 35593.6989427372
```

```
In [ ]: def perform_regression(x, y):
X_train, X_test, y_train, y_test= train_test_split(X, y, test_size=0.2)
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
RMSE = np.sqrt(mean_squared_error(y_test, y_preds))
return RMSE, y_preds
```

```
Out [ ]: for clf in [RidgeCV(), ElasticNetCV(), LassoCV()]:
print(clf.__class__.__name__)
RMSE, y_preds = perform_regression(X, y)
print(RMSE)
```

```
Out [ ]: RidgeCV
56373.137803244654
ElasticNetCV
56645.6301979143
LassoCV
44998.85202724976
```

```
In [ ]: RF_params= {'n_estimators': [100, 250, 500], 'max_depth': [2, 5, 10, None]}
```

```
reg = RandomForestRegressor()
grid = GridSearchCV(reg, RF_params, cv=5, scoring = make_scorer(mean_squared_error))
grid.fit(X, y)
```

```
Out [ ]: GridSearchCV(cv=5, error_score=nan,
estimator=RandomForestRegressor(bootstrap=True, ccp_alpha=0.0,
criterion='gini', max_depth=None, max_features='auto',
max_leaf_nodes=None, max_samples=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=100,
n_jobs=None, oob_score=False, random_state=42, verbose=0,
warm_start=False),
id=1,
param_grid={'max_depth': [2, 5, 10, None]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
scoring=make_scorer(mean_squared_error), verbose=0)
```

```
In [ ]: print(grid.best_estimator_)
print(grid.best_estimator_.get_params())
print(grid.best_score_)
print(grid.best_estimator_.get_params())
```

```
Out [ ]: ({"max_depth": 2, "n_estimators": 100},
RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
max_depth=2, max_features='auto', max_leaf_nodes=None,
max_samples=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=1,
min_samples_split=2, min_weight_fraction_leaf=0.0,
n_estimators=100, n_jobs=None, oob_score=False,
random_state=None, verbose=0, warm_start=False))
```

```
Out [ ]: 0.16589498 7256362
```

```
In [ ]: grid.best_estimator_
```

```
Out [ ]: make_scorer(mean_squared_error)
```

```
In [ ]: y_pred = grid.best_estimator_.predict(X)
```

```
In [ ]: GB_params= {'n_estimators': [100, 250, 500], 'max_depth': [2, 3, 5], 'learning_rate': [0.1, 0.25, 0.5]}
```

```
gb = GradientBoostingRegressor()
grid = GridSearchCV(gb, GB_params, cv=5, scoring = make_scorer(mean_squared_error))
grid.fit(X, y)
```

```
Out [ ]: GridSearchCV(cv=5, error_score=nan,
estimator=GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0,
criterion='friedman_mse', max_depth=3,
init=None, learning_rate=0.5, loss='ls', max_depth=3,
max_features='auto', 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, n_estimators=500,
n_iter_no_change=None, presort='deprecated',
random_state=None, subsample=1.0, tol=0.0001,
validation_fraction=0.1, verbose=0, warm_start=False),
id=1,
param_grid={'learning_rate': [0.1, 0.25, 0.5],
'n_estimators': [100, 250, 500]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
scoring=make_scorer(mean_squared_error), verbose=0)
```

```
In [ ]: print(grid.best_params_)
print(grid.best_estimator_)
print(grid.best_index_)
print(grid.best_score_)
print(np.sqrt(grid.best_score_))
```

```
Out [ ]: ({"learning_rate": 0.5, "max_depth": 3, "n_estimators": 500},
GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0, criterion='friedman_mse',
max_depth=3, max_features='auto', max_leaf_nodes=None,
max_samples=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=1,
min_samples_split=2, min_weight_fraction_leaf=0.0,
n_estimators=500, n_iter_no_change=None, presort='deprecated',
random_state=None, subsample=1.0, tol=0.0001,
validation_fraction=0.1, verbose=0, warm_start=False))
```

```
Out [ ]: 23
1344895041.1345449
45908.40817658669 38915.37795183987
```

```
In [ ]: y_prediction = grid.best_estimator_.predict(X)
```

```
In [ ]: print(XT.shape)
print(housing_test.shape)
print(y_pred.shape)
print(y_prediction.shape)
```

```
Out [ ]: (1459, 80)
(1459, 80)
(1459,)
```

```
In [ ]: pd.DataFrame(y_prediction)
pd.DataFrame(y_pred)
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 = pd.DataFrame(y_prediction)
GridBoost = pd.concat([Identification, SalesPrice0], axis=1)
RandomForest
```

```
Out [ ]: Id      0
SalePrice 133664.130163
1 140249.805542
2 140296.162097
3 150991.532061
4 266422.704043
...
1459 211642.669999
1459 211642.669999
```

```
In [ ]: Identification = housing_test['Id'].copy()
SalesPrice0 = pd.DataFrame(y_pred)
SalesPrice1 = pd.DataFrame(y_prediction)
RandomForest = pd.concat([Identification, SalesPrice0], axis=1)
SalesPrice0 = pd.DataFrame(y_prediction)
SalesPrice1 =
```