



Capstone Project

Home Sales Prices

By: Ryan Mark

Contents



Executive Summary



Exploratory Data Analysis (EDA)



Design and modeling methods



Results



Recommendations

Executive Summary

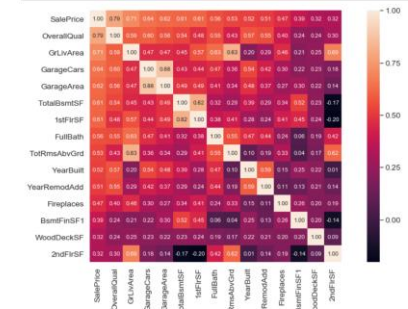
- The Ames, Iowa housing dataset captures various details about residential properties in that area from 2006 to 2010. The objective is to facilitate the use of machine learning models to predict the value of home sales prices.

Exploratory Data Analysis (EDA)

- The initial Exploratory Data Analysis (EDA) was focused on understanding the target variable 'SalePrice'.
- From this information it can be determined that the volume of prices was between 100k and 200k.
- This output indicated that Overall Quality had a strong positive correlation with Sale Price.

```
#SalePrice correlation matrix
k = 15 #number of variables for heatmap
plt.figure(figsize=(15,8))
corrmat = train_df.corr()

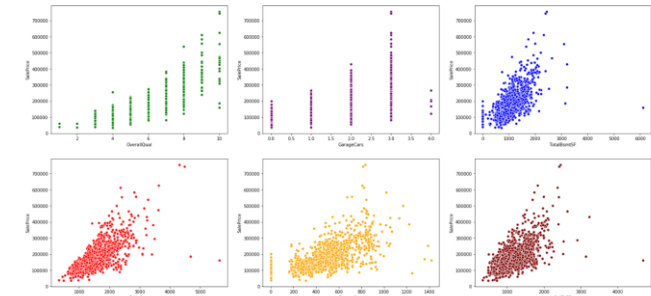
# picking the top 25 correlated features
cols = corrmat.nlargest(k, 'SalePrice')['SalePrice'].index
cm = np.corrcoef(train_df[cols].values.T)
sns.set(font_scale=1.5)
hm = sns.heatmap(cm, cbar=True, annot=True, square=True, ftext='.2f', annot_kws={'size': 10}, yticklabels=cols.values, xticklabels=cols.values)
plt.show()
```



```
# Top 6 correlated columns to SalePrice graphed
fig, axes = plt.subplots(2, 3, figsize=(25, 12), sharex=False)

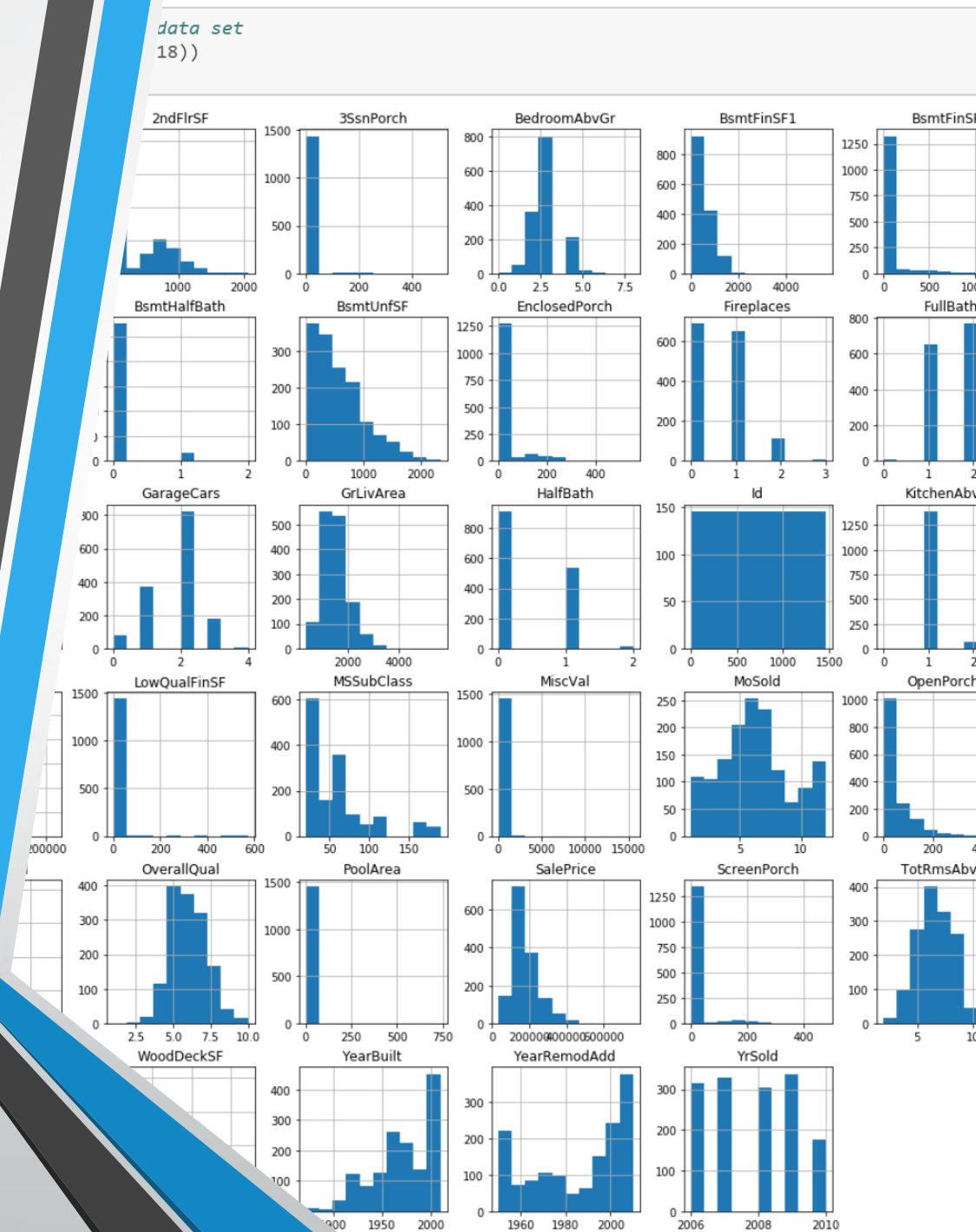
sns.scatterplot(x='OverallQual', y='SalePrice', data=train_df, color='green', ax=axes[0, 0])
sns.scatterplot(x='GrLivArea', y='SalePrice', data=train_df, color='red', ax=axes[1, 0])
sns.scatterplot(x='GarageCars', y='SalePrice', data=train_df, color='purple', ax=axes[0, 1])
sns.scatterplot(x='GarageArea', y='SalePrice', data=train_df, color='orange', ax=axes[1, 1])
sns.scatterplot(x='TotalBsmSF', y='SalePrice', data=train_df, color='blue', ax=axes[0, 2])
sns.scatterplot(x='1stFlrSF', y='SalePrice', data=train_df, color='maroon', ax=axes[1, 2])

plt.show()
```



Review research design and modeling methods

- The modeling methods used in our analysis consisted of Linear Regression, Ridge Regression, Random Forest, and Gradient Boosting. The intent was to understand the performance of Random Forest and Gradient Boosting in respect to the previously used models of Linear / Ridge regression. Through Random Forest Regressor a mean prediction can be obtained from the multiple decision trees within that model. Through Gradient Boosting, it can strengthen a model with weak predictions and make it better.



Results

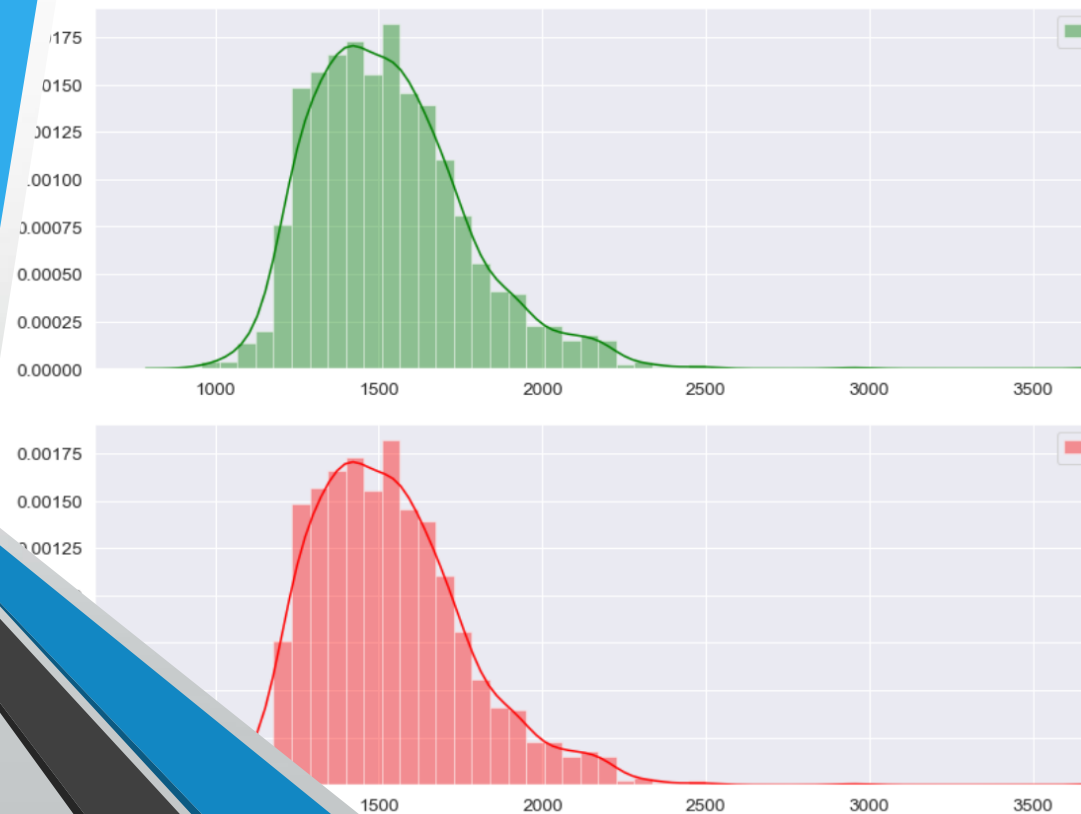
- Robust Scaling was used on all models to scale the data based on the interquartile range. This was done to address the robust approach of this scaling method to outliers. This scaling approach was used on the Random Forest and Gradient Boosting models initially as a test. It should be known that in principle tree models do not require scaling it is unaffected by monotonic transformations.
- Our results indicated that the Gradient Boosting Model performed the best in terms of RMSE. This was further evidenced by the Kaggle score submission which was the highest score achieved (score = .164680). This indicates that the predictor values in our top model was the best in terms of having the lowest prediction error. Furthermore, Both the Random Forest and Gradient Boosting performed distinctively better than the Linear Regression and Ridge Regression models.

```
overlap
figsize=(15, 5))
t(lr_final['SalePrice'], label='Predicted', axlabel = False, color='green')
)

t overlap
:(figsize=(15, 5))
lot(rr_final['SalePrice'], label='Predicted', axlabel = False, color='red')
d()
.)

lot overlap
ire(figsize=(15, 5))
tplot(rf_final['SalePrice'], label='Predicted', axlabel = False, color='purple')
end()
)w()

plot overlap
.gure(figsize=(15, 5))
istplot(gbr_final['SalePrice'], label='Predicted', axlabel = False, color='blue')
legend()
show()
```



Exposition, problem description, and management recommendations

- The comparison of Linear Regression, Ridge Regression, Random Forest, and Gradient Boosting indicated a significant performance when utilizing multiple tree methods. In particular it was observed that Random Forest and Gradient Boosting superior predictive results compared to our base regression models. Based on the model comparisons, Gradient Boosting is the best model for predicting home prices.

```
data = {'Model':      ['Linear Regression', 'Ridge Regression', 'Random Forest_2', 'Gradient Boosting_2'],  
        'RMSE Score': [.454, .454, .378, .354, 31341.180, 29687.79],  
        'Kaggle Score': [4.70255, 4.70256, 10.76463, 11.21304, .16601, .16468],  
        'Scaling?':    ['Scaled', 'Scaled', 'Scaled', 'Scaled', 'Not Scaled', 'Not Scaled']  
    }
```

```
results = pd.DataFrame(data, columns = ['Model', 'RMSE Score', 'Kaggle Score', 'Scaling?'])  
results.style.highlight_max(color='red').highlight_min(color='green')
```

	Model	RMSE Score	Kaggle Score	Scaling?
0	Linear Regression	0.454000	4.702550	Scaled
1	Ridge Regression	0.454000	4.702560	Scaled
2	Random Forest	0.378000	10.764630	Scaled
3	Gradient Boosting	0.354000	11.213040	Scaled
4	Random Forest_2	31341.180000	0.166010	Not Scaled
5	Gradient Boosting_2	29687.793000	0.164680	Not Scaled