

	<ul> <li>Regression Models</li> <li>For regression modelling, we will be exploring the following models:</li> <li>1. Linear Regression (simple)</li> <li>2. Ridge Regression - Penalized for the sum of squared value of the weights (L2)</li> <li>3. Lasso Regression - Penalized for the sum of absolute values of the weights (L1)</li> <li>4. SVR Regression - Best fit hyperplane</li> <li>5. Random Forest Regression - Decision Tree</li> </ul>
In [27]:	Using GridSearchCV on the train data, we are able to optimize each model via their hyper parameters.  Notice for the pipeline, all models will follow the same pre-processor as illustrated in the models.  1. Linear Regression  from math import sqrt  #1. Model Creation model = LinearRegression()
	<pre>pipeline = Pipeline(     steps=[("preprocessor", preprocessor),         ('model', model)] )  #2. Pre-Processing Parameters to Hyperparameter Search param_grid = {     #Simple linear regression no need Gridsearch }  best_modellr = GridSearchOp(pipeline,param_grid,cv) yhat_train = best_modellr.predict(X_train)</pre>
	<pre>yhat_test = best_modellr.predict(X_test)  scoreEval(y_train,yhat_train,y_test,yhat_test)  set_config(display="diagram") pipeline # click on the diagram below to see the details of each step  Fitting 30 folds for each of 1 candidates, totalling 30 fits Best params: {} R2 Score for train: 0.703 R2 Score for test: 0.686 RMSE Score for train: 200558.25</pre>
Out[27]:	<pre>preprocessor: ColumnTransformer     num</pre>
In [28]:	<pre>2. Ridge Regression  from sklearn.linear_model import Ridge  # fit a linear regression model to the data model = Ridge()  pipeline = Pipeline(     steps=[("preprocessor", preprocessor),</pre>
Out[28]:	<pre>yhat_train = best_modelridge.predict(X_train) yhat_test = best_modelridge.predict(X_test)  scoreEval(y_train, yhat_train, y_test, yhat_test)  set_config(display="diagram") pipeline # click on the diagram below to see the details of each step  Fitting 30 folds for each of 8 candidates, totalling 240 fits Best params: {'model_solver': 'sparse_cg'} R2 Score for train: 0.703 R2 Score for test: 0.685 RMSE Score for test: 203890.18</pre>
In [29]:	* SimpleImputer
	<pre>steps=[("preprocessor", preprocessor),     ('model', model)]  #2. Pre-Processing Parameters to Hyperparameter Search param_grid = {     "modelselection": ['cyclic', 'random'], #Simple linear regression no need Gridse }  best_modellasso = GridSearchOp(pipeline,param_grid,cv) yhat_train = best_modellasso.predict(X_train) yhat_test = best_modellasso.predict(X_test)  scoreEval(y_train,yhat_train,y_test,yhat_test)  set_config(display="diagram") pipeline # click on the diagram below to see the details of each step  Fitting 30 folds for each of 2 candidates, totalling 60 fits Best params:</pre>
Out[29]:	<pre>{'modelselection': 'cyclic'} R2 Score for train: 0.703 R2 Score for test: 0.686 RMSE Score for train: 200558.25 RMSE Score for test: 203841.56  Pipeline  Preprocessor: ColumnTransformer  num</pre>
In [30]:	<pre>from sklearn.svm import LinearSVR  # fit a linear regression model to the data model = LinearSVR()  pipeline = Pipeline(     steps=[("preprocessor", preprocessor),         ('model', model)] )  #2. Pre-Processing Parameters to Hyperparameter Search param_grid = {     #Linear SVC no need Gridsearch }  best_modellsvr = GridSearchOp(pipeline,param_grid,cv) yhat_train = best_modellsvr.predict(X_train) yhat_test = best_modellsvr.predict(X_test)  scoreEval(y_train,yhat_train,y_test,yhat_test)  set_config(display="diagram") pipeline # click on the diagram below to see the details of each step  Fitting 30 folds for each of 1 candidates, totalling 30 fits Best params: {}</pre>
Out[30]: In [31]:	R2 Score for train: -1.480 R2 Score for test: -1.520 RMSE Score for train: 579548.01 RMSE Score for test: 577071.52  Pipeline  preprocessor: ColumnTransformer  num  cat  remainder  SimpleImputer  OneHotEncoder  passthrough  LinearSVR  from sklearn.linear model import SGDRegressor
	<pre># fit a linear regression model to the data model = SGDRegressor()  pipeline = Pipeline(     steps=[("preprocessor", preprocessor),         ('model', model)] )  #2. Pre-Processing Parameters to Hyperparameter Search param_grid = {     "model_loss": ['squared_error', 'huber', 'squared_epsilon_insensitive'], #Simple 1.     "model_ penalty": ['12','11','elasticnet'],     "model_ alpha": [0.0001,0.001] }  best_modelsgd = GridSearchOp(pipeline,param_grid,cv) yhat_train = best_modelsgd.predict(X_train) yhat_test = best_modelsgd.predict(X_test)</pre>
Out[31]:	scoreEval(y_train,yhat_train,y_test,yhat_test)  set_config(display="diagram") pipeline # click on the diagram below to see the details of each step  Fitting 30 folds for each of 18 candidates, totalling 540 fits  Best params: {'model_alpha': 0.001, 'model_loss': 'squared_error', 'model_penalty': '12'}  R2 Score for train: 0.693  R2 Score for test: 0.666  RMSE Score for train: 203996.65  RMSE Score for test: 210228.45  Pipeline  Pipeline  Preprocessor: ColumnTransformer  Passthrough  RobustScaler  SGDRegressor  PSGDRegressor
In [32]:	<pre>from sklearn.ensemble import RandomForestRegressor  # fit a linear regression model to the data model = RandomForestRegressor(random_state=42)  pipeline = Pipeline(     steps=[("preprocessor", preprocessor),</pre>
Out[32]:	yhat_test = best_modelrfr.predict(X_test)  scoreEval(y_train, yhat_train, y_test, yhat_test)  set_config(display="diagram") pipeline # click on the diagram below to see the details of each step  Fitting 30 folds for each of 12 candidates, totalling 360 fits  Best params: {'model_max_depth': 6, 'model_n_estimators': 100}  R2 Score for train: 0.844  R2 Score for test: 0.817  RMSE Score for train: 145192.10  RMSE Score for test: 155511.57
In [33]:	* RandomForestRegressor  6. Model Evaluation Performance  modelList = ['Linear', 'Ridge', 'Lasso', 'SVR', 'SGD', 'RF'] #Storage for model performance  ind = np.arange(len(modelList)) # the x locations for the groups width = 0.35 # the width of the bars  fig, ax = plt.subplots(figsize=(10,5)) rects1 = plt.bar(ind, train r2, width, color='b')
Out[33]:	rects1 = plt.bar(ind, train_r2, wldth, color='b') rects2 = plt.bar(ind + width / 2) ax.set_xticks(ind + width / 2) ax.set_xticklabels((modelList))  plt.hlines(max(train_r2), 0, len(modelList), colors="g", linestyles='dotted') ax.text(len(modelList), max(train_r2), 'Max. Train_R2:(:.2f)'.format(max(train_r2)), : plt.hlines(max(test_r2), 0, len(modelList), colors="b", linestyles='dotted') ax.text(len(modelList), max(test_r2), 'Max. Test_R2:(:.2f)'.format(max(test_r2)), fon(ax.legend((rects1[0], rects2[0]), ('Train_R2', 'Test_R2'),loc='lower right') plt.title("Regression Model Comparison-Test Data") plt.xlabel("Models") plt.ylabel("Scores")  Regression Model Comparison-Test Data  Regression Model Comparison-Test Data  Regression Model Comparison-Test Data
In [34]:	Feature Importance: Random Forest Regressor  featModel = best_modelrfr.best_estimatornamed_steps["model"] importance = featModel.feature_importances_ print(importance) print(len(importance)) combination = pd.Series(importance, columns) plt.figure(figsize=(10,10)) combination.sort_values().plot.barh(color='red') plt.title("Visualized Feature Importance: {}".format('Random Forest Regressor'))  [3.08446947e-04 2.84292772e-03 3.09702908e-01 4.85068029e-03
Out[34]:	3.11077132e-04 8.68110960e-03 2.37645782e-04 3.90256830e-01 1.71221389e-03 1.80655868e-02 7.20589804e-04 4.97605313e-03 1.62640121e-01 6.01454165e-02 2.38659151e-04 4.18568467e-04 1.59003708e-02 1.79907951e-02] 18  Text(0.5, 1.0, 'Visualized Feature Importance: Random Forest Regressor')  Visualized Feature Importance: Random Forest Regressor  num_grade  num_sqft_living  num_long  num_vr_built  cat_x1_1  cat_x1_0
	num_zipcode -  num_sqft_lot -  num_bathrooms -  num_sqft_basement -  num_yr_renovated -  cat_x0_Renovated -  num_floors -  num_bedrooms -  cat_x0_Not Renovated -  num_condition -
In [35]:	#Plotting the relationship between prices and grade  plt.figure(figsize=(6.5,5)) sns.boxplot(x='grade', y='price', data=housedf) plt.xlabel("Grade", fontsize=15) plt.ylabel("Price (\$)", fontsize=15) plt.xticks(fontsize=13) plt.yticks(fontsize=13) plt.legend() plt.title("Box plot of different grades(assume categorized)")  plt.gca().spines['right'].set_visible(False)  No handles with labels found to put in legend.
	Box plot of different grades(assume categorized)  7 - 6 -
	(\$\frac{1}{3}\) 2 1 1 1 3 4 5 6 7 8 9 10 11 12 13  Grade
In [36]:	
In [36]:	Baseline-Score (Reference)    mean_price = y.mean()   print("mean_price={}".format(mean_price))   ystupid = np.full(len(y),mean_price)   ystupid   print("R2 Score for stupid baseline is={}".format(r2_score(y,ystupid)))   mean_price=540088.1417665294
In [36]:	Baseline-Score (Reference)    mean_price = y.mean()   print("mean_price")   ystupid = np.full(len(y), mean_price)   ystupid   print("R2 Score for stupid baseline is={}".format(r2_score(y, ystupid)))     mean_price=540088.1417665294   R2 Score for stupid baseline is=0.0  7. Model Selection: Random Forest Regressor Model   Back to Main
In [36]:	Baseline-Score (Reference)  mean_price = y.mean() print("mean price={}".format(mean_price)) ystupid = np.full(len(y),mean_price) ystupid = np.full(len(y),mean_price) ystupid = ror stupid baseline is={}".format(r2_score(y,ystupid)))  mean_price=540088.1417665294 R2_Score_for_stupid baseline is=0.0  7. Model Selection: Random Forest Regressor Model  Back to Main  • For the modelling pipeline, the initial approach is to bin the renovation year into categories and yet, to retain the numeric relationship for the houses which have renovated. This will provide more information to the model rather than just renovated/not renovated. • To remove the affects of 0 year in year renovated, I have replaced it with the average of the column so that during standardscaling, these rows of data will not impact the feature bias. • Based on the model performance and RMSE, we can see clearly that the best model is the random forest regressor model. • Interestingly, from the feature importance, despite all the feature engineering to the yr_renovated, it is seen that the derivatives from the feature variable are considered low in importance for prediction of the prices. Instead the grade of the houses are most important. • From the boxplot, it seemed like the grades should be treated as ordinal categories and have been picked up well by the random forest regressor as seen by significant differences in prices with higher grades.
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
	Baseline-Score (Reference)    Definit ("mean price = y .mean()
In [36]:	Baseline-Score (Reference)    Definit ("mean price = y .mean()
	Baseline-Score (Reference)    Definit ("mean price = y .mean()