

IBM MACHINE LEARNING

SUPERVISED MACHINE LEARNING: REGRESSION MODELS FOR HOUSE PRICE PREDICTION



Fatima, Sayeda

9/24/2022

Table of Contents

1) Project Overview	2
2) About the Dataset	3
2a) Brief description of the data set you chose:.....	3
2b) Summary of Data Attributes	3
3) Main Objectives of Analysis	4
4) Data Exploration, Data Cleansing and Features Engineering	5
4a) Data Exploration	5
4b) Data Cleansing & Features Engineering.....	9
5) Summary of Training Different Regression Models	16
5a) Machine Learning Regression Algorithm Development Approach.....	16
5b) Summarizing Employed Models	20
1) Ridge Regression (RR) Models.....	20
2) Lasso Regression (LR) Models	23
3) Elastic-Net (EN) Models.....	26
4) Extreme Gradient Boosting (XGB) Regression Models.....	29
6) Key Findings to the Main Objectives of Analysis.....	32
6a) Result Summary	32
6b) Recommended Model and Justification	32
6c) Summarizing Model Drivers	33
6d) Enlisting Top Contributory Factors	33
6e) Visualizing Top Contributory Factors Driving Sale Value.....	34
6f) Sale Price Prediction on Test/New Data.....	35
7) Future Directions	36
7a) Possible Flaws in Chosen Model	36
7b) Recommendations	36
8) Useful Links	37
8a) Link to Other Useful Models	37
8b) Github Link to Assignment Notebook and Other Files	37
References.....	38

1) Project Overview

A fundamental issue real estate business face is assessing sale price based on house attributes. Grounded on hedonic price modelling theory, not only do neighbourhood-specific characteristics but also unit-specific attributes greatly drive house prices (Herath and Maier, 2010). Hence, identification of key factors driving sale prices of real estate is crucial to facilitate informed purchase decisions. It is here that Regression Machine Learning models can be very useful to gain deeper insight into underlying factors as well as their relationship in driving and estimating fair sale price of houses.

Hence, the main aim of the following regression modelling and analysis approach is to enable the business to:

- * Find best regression model for house price prediction
- * Identify different factors influencing house prices
- * Predict sale price based on contributory factors

2) About the Dataset

2a) Brief description of the data set you chose:

This project uses a hypothetical dataset 'Ames, Iowa Housing Dataset' which was downloaded from the following link:

<https://www.kaggle.com/prevek18/ames-housing-dataset>

2b) Summary of Data Attributes

The dataset exhibits 2,930 data points (rows) and 82 features (columns) reflecting on housing characteristics.

The data also comes with 'SalePrice' Column which represents the Class requiring prediction.

#	Column	Non-Null Count	Dtype
0	Order	2930 non-null	int64
1	PID	2930 non-null	int64
2	MS SubClass	2930 non-null	int64
3	MS Zoning	2930 non-null	object
4	Lot Frontage	2440 non-null	float64
5	Lot Area	2930 non-null	int64
6	Street	2930 non-null	object
7	Alley	198 non-null	object
8	Lot Shape	2930 non-null	object
9	Land Contour	2930 non-null	object
10	Utilities	2930 non-null	object
11	Lot Config	2930 non-null	object
12	Land Slope	2930 non-null	object
13	Neighborhood	2930 non-null	object
14	Condition 1	2930 non-null	object
15	Condition 2	2930 non-null	object
16	Bldg Type	2930 non-null	object
17	House Style	2930 non-null	object
18	Overall Qual	2930 non-null	int64
19	Overall Cond	2930 non-null	int64
20	Year Built	2930 non-null	int64
21	Year Remod/Add	2930 non-null	int64
22	Roof Style	2930 non-null	object
23	Roof Matl	2930 non-null	object
24	Exterior 1st	2930 non-null	object
25	Exterior 2nd	2930 non-null	object
26	Mas Vnr Type	2907 non-null	object
27	Mas Vnr Area	2907 non-null	float64
28	Exter Qual	2930 non-null	object
29	Exter Cond	2930 non-null	object
30	Foundation	2930 non-null	object
31	Bsmt Qual	2850 non-null	object
32	Bsmt Cond	2850 non-null	object
33	Bsmt Exposure	2847 non-null	object
34	BsmtFin Type 1	2850 non-null	object
35	BsmtFin SF 1	2929 non-null	float64
36	BsmtFin Type 2	2849 non-null	object
37	BsmtFin SF 2	2929 non-null	float64
38	Bsmt Unf SF	2929 non-null	float64
39	Total Bsmt SF	2929 non-null	float64
40	Heating	2930 non-null	object
41	Heating QC	2930 non-null	object
42	Central Air	2930 non-null	object
43	Electrical	2929 non-null	object

3) Main Objectives of Analysis

Real estate business performance is largely dependent on paying fair price of assets to prevent overpriced purchases and minimize loss. Hence, these businesses are continuously faced with the challenge to estimate realistic house prices and often rely on manual application of Hedonic Price Method (HPM) or hedonic regression analysis. Consequently, driven by HPM, this analysis is targeted towards answering the following queries

- What are the various contributory factors which drive house prices in a given area?
- Based on important factors, what will be projected price for different housing units?

As a consequence, implementation of an automated machine learning (ML) HPM regression modelling process will enable the organization to:

- identify key underlying factors which can appreciate or depreciate property values
- save valuable resources and funds in purchasing properties at right values
- effortlessly employ best ML HPM model and generate report with the click of a button

4) Data Exploration, Data Cleansing and Features Engineering

Since the quality of any machine learning model highly depends on quality of data, hence, this stage is not only most important but is also time consuming. Hence, it was conducted in a step-by-step manner.

4a) Data Exploration

- Data was first loaded into pandas dataframe

Section A: About the Dataset

1a) Brief description of the data set you chose:

This project uses dataset for 'Ames Housing' which may also be downloaded from the following link:

<https://www.kaggle.com/datasets/prevek18/ames-housing-dataset>

This dataset represents housing data for Ames, Iowa region in between 2006 to 2010.

1b) Summary of Data Attributes:

The dataset exhibits 2930 data points (rows) and the following 82 features (columns):

```
1 raw_data = pd.read_csv('C:/Users/fatima.s/Documents/PythonScripts/DATA SCIENCE/IBM Machine Learning Intermediate/MODULE 2 SU
2 raw_data
```

	Order	PID	MS SubClass	MS Zoning	Lot Frontage	Lot Area	Street	Alley	Lot Shape	Land Contour	...	Pool Area	Pool QC	Fence	Misc Feature	Misc Val	Mo Sold	Yr Sold	Sale Type	Sa
0	1	526301100	20	RL	141.0	31770	Pave	NaN	IR1	Lvl	...	0	NaN	NaN	NaN	0	5	2010	WD	Norm
1	2	526350040	20	RH	80.0	11622	Pave	NaN	Reg	Lvl	...	0	NaN	MnPrv	NaN	0	6	2010	WD	Norm
2	3	526351010	20	RL	81.0	14267	Pave	NaN	IR1	Lvl	...	0	NaN	NaN	Gar2	12500	6	2010	WD	Norm
3	4	526353030	20	RL	93.0	11160	Pave	NaN	Reg	Lvl	...	0	NaN	NaN	NaN	0	4	2010	WD	Norm
4	5	527105010	60	RL	74.0	13830	Pave	NaN	IR1	Lvl	...	0	NaN	MnPrv	NaN	0	3	2010	WD	Norm
...
2925	2926	923275080	80	RL	37.0	7937	Pave	NaN	IR1	Lvl	...	0	NaN	GdPrv	NaN	0	3	2006	WD	Norm
2926	2927	923276100	20	RL	NaN	8885	Pave	NaN	IR1	Low	...	0	NaN	MnPrv	NaN	0	6	2006	WD	Norm
2927	2928	923400125	85	RL	62.0	10441	Pave	NaN	Reg	Lvl	...	0	NaN	MnPrv	Shed	700	7	2006	WD	Norm
2928	2929	924100070	20	RL	77.0	10010	Pave	NaN	Reg	Lvl	...	0	NaN	NaN	NaN	0	4	2006	WD	Norm
2929	2930	924151050	60	RL	74.0	9627	Pave	NaN	Reg	Lvl	...	0	NaN	NaN	NaN	0	11	2006	WD	Norm

2930 rows x 82 columns

- A method was created to conduct preliminary analysis including computation of:
 - ☞ Descriptive statistics to summarize shape of a dataset's distribution, its dispersion and central tendency
 - ☞ Data analysis to depict data types, skewness, kurtosis, etc to facilitate subsequent data cleansing

```

1 def analysis(*name): # This Method will extract dataframe by name
2     n = name # Extract Dataframe by Name...this will create a 3d tuple
3     n = (n[0]) # Convert Tuple to To Dataframe
4     df_name = [x for x in globals() if globals()[x] is n][0] # Extract Name of Imported Dataframe to print Later
5
6     # Perform Statistics
7     stats = n.describe(include = 'all').transpose()
8     stats = stats.fillna(0) # Replace all Nan Values with Zero
9
10    # Data Analysis
11    obs = n.shape[0]
12    types = n.dtypes
13    counts = n.apply(lambda x: x.count())
14    distincts = n.apply(lambda x: x.unique().shape[0])
15    nulls = n.apply(lambda x: x.isnull().sum())
16    uniques = n.apply(lambda x: [x.unique()])
17    per_nulls = (n.isnull().sum()/ obs) * 100
18    skewness = n.skew()
19    kurtosis = n.kurt()
20    corr = n.corrwith(n[target])# "SalePrice"
21    #corr = corr.to_string()
22
23    # Transform Data Analysis to Dataframe
24    analyze = pd.DataFrame(columns=['Columns', 'types', 'counts', 'distincts', 'nulls', '% nulls', 'uniques', 'skewness', 'ku
25    analyze['types'] = types
26    analyze['counts'] = counts
27    analyze['distincts'] = distincts
28    analyze['nulls'] = nulls
29    analyze['% nulls'] = per_nulls
30    analyze['uniques'] = uniques
31    analyze['skewness'] = skewness
32    analyze['kurtosis'] = kurtosis
33    analyze['Corr_Sales'] = corr
34    analyze['Columns'] = analyze.index
35    analyze = analyze.fillna(0).sort_values(by=['Corr_Sales', 'skewness'], ascending=False) #Fill Remaining Missing Values wi
36    analyze = analyze.replace(["NaN"], 0).sort_values(by='Corr_Sales', ascending=False)
37    analyze = analyze.reset_index(drop=True)
38
39    print(colored("\nData Analysis for: ", 'green', attrs=['bold'])
40          +colored(df_name, 'red', attrs=['bold'])
41          + colored("\nData Shape:", 'green', attrs=['bold'])
42          +colored(obs, 'magenta', attrs=['bold'])
43          )
44
45    return analyze, stats

```

Data Analysis for: raw_data
Data Shape: 2930

Summary Statistics

	count	unique	top	freq	mean	std	min	25%	50%	75%	max
Order	2930.0	0	0	0	1.465500e+03	8.459625e+02	1.0	7.332500e+02	1465.5	2.197750e+03	2.930000e+03
PID	2930.0	0	0	0	7.144645e+08	1.887308e+08	526301100.0	5.284770e+08	535453620.0	9.071811e+08	1.007100e+09
MS SubClass	2930.0	0	0	0	5.738737e+01	4.263802e+01	20.0	2.000000e+01	50.0	7.000000e+01	1.900000e+02
MS Zoning	2930.0	7	RL	2273	0.000000e+00	0.000000e+00	0.0	0.000000e+00	0.0	0.000000e+00	0.000000e+00
Lot Frontage	2440.0	0	0	0	6.922459e+01	2.336533e+01	21.0	5.800000e+01	68.0	8.000000e+01	3.130000e+02
...
Mo Sold	2930.0	0	0	0	6.216041e+00	2.714492e+00	1.0	4.000000e+00	6.0	8.000000e+00	1.200000e+01
Yr Sold	2930.0	0	0	0	2.007790e+03	1.316613e+00	2006.0	2.007000e+03	2006.0	2.009000e+03	2.010000e+03
Sale Type	2930.0	10	WD	2536	0.000000e+00	0.000000e+00	0.0	0.000000e+00	0.0	0.000000e+00	0.000000e+00
Sale Condition	2930.0	6	Normal	2413	0.000000e+00	0.000000e+00	0.0	0.000000e+00	0.0	0.000000e+00	0.000000e+00
SalePrice	2930.0	0	0	0	1.807961e+05	7.988669e+04	12789.0	1.295000e+05	160000.0	2.135000e+05	7.550000e+05

82 rows × 11 columns

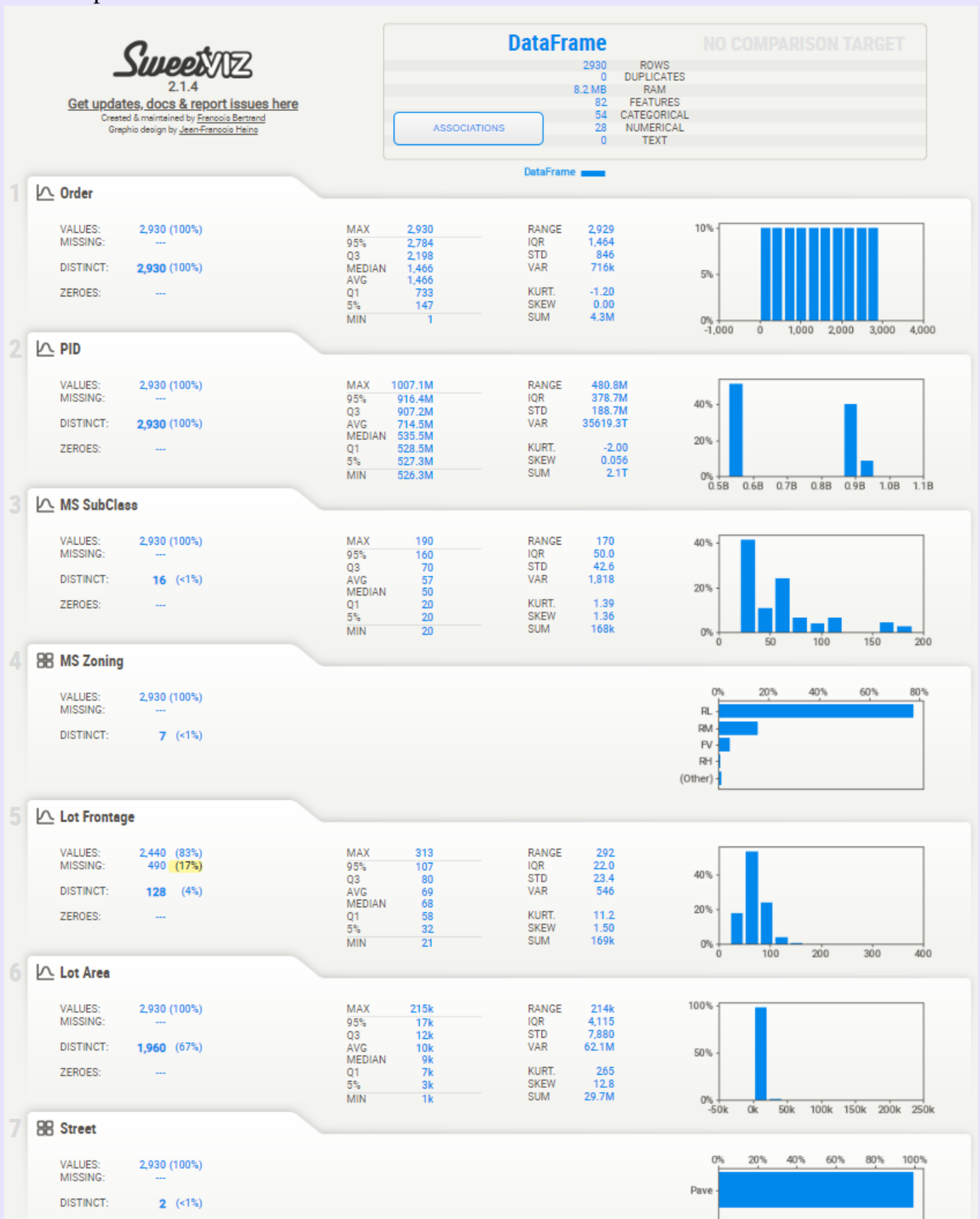
```
1 print(colored("\nData Analysis Summary:", 'cyan', attrs=['bold']))
2 analyze
```

Data Analysis Summary:

	Columns	types	counts	distincts	nulls	% nulls	uniques	skewness	kurtosis	Corr_Sales
0	SalePrice	int64	2930	1032	0	0.00000	0	1.743500	5.118900	1.000000
1	Overall Qual	int64	2930	10	0	0.00000	0	0.190634	0.052412	0.799262
2	Gr Liv Area	int64	2930	1292	0	0.00000	0	1.274110	4.137838	0.706780
3	Garage Cars	float64	2929	7	1	0.03413	0	-0.219836	0.244969	0.647877
4	Garage Area	float64	2929	604	1	0.03413	0	0.241994	0.951023	0.640401
...
77	MS SubClass	int64	2930	16	0	0.00000	0	1.357579	1.386775	-0.085092
78	Overall Cond	int64	2930	9	0	0.00000	0	0.574429	1.491450	-0.101697
79	Kitchen AbvGr	int64	2930	4	0	0.00000	0	4.313825	19.869743	-0.119814
80	Enclosed Porch	int64	2930	183	0	0.00000	0	4.014446	28.487205	-0.128787
81	PID	int64	2930	2930	0	0.00000	0	0.055886	-1.995146	-0.246521

82 rows × 10 columns

- Additional Automated Exploratory Data Analysis was performed using Sweetviz to realize visual representation.



4b) Data Cleansing & Features Engineering

In machine learning, feature selection is the method to reduce the number of input variables during developing predictive modelling. This reduction in input variables is necessary not only to minimize computational cost of modelling but also to achieve improved performance of the model itself.

Among widely practices feature selection approaches include statistical-based feature selection methods which use statistical measures to evaluate relationship between each input variable and the target variable and then select those exhibiting strongest relationship with the latter. While these methods can be both speedy and effective, however, the ultimate choice of statistical measure is largely dependent on data types of both of these variables.

Irrespective of the statistical measure being employed, two dominant feature selection techniques, that is supervised and unsupervised, exist where the former can be further categorized into wrapper, filter and intrinsic techniques. Filter-based feature selection methods employs statistical measures to evaluate correlation between input and output variables so that those exhibiting highest correlations are selected. Statistical measures employed in filter-based feature selection are normally univariate in nature since they evaluate relationship of single input variables one by one with target variable, disregarding their interaction with each other.

Consequently, adopting filter-based feature selection methods, the housing price prediction model approached filter engineering in three steps.

- 1) Data Encoding
- 2) Managing Multicollinearity
- 3) Final Data Cleansing
- 4) Applying Outlier Treatment

- 1) Prior to final features selection, data encoding of object or string columns was carried out to facilitate any statistical computation during features selection process. Hence, after copying original dataset, an automated method was created and employed to encode object data using Scikit-learn label encoder.

Data Cleansing Actions

Method to Encode Object Type Columns :
 1) List Object Type Columns & Encode Data
 2) Make Decoder to Decode Encoded Data

```

1 # Method to encode object/string columns
2 def encoder(*name):
3     # Accept an argument, return a value.
4     n = name # Extract Dataframe by Name...this will create a 3d tuple
5     n = (n[0]) # Convert Tuple to To Dataframe
6     df_name = [x for x in globals() if globals()[x] is n][0] # Extract Name of Imported Dataframe
7
8     # 1) List all Object/String Columns
9     from sklearn import preprocessing
10    cat_columns = n.select_dtypes(include=[object]) # Get Object Type Columns to Convert to Encoded Categories
11
12    categorical_column = list(cat_columns.columns) # List of columns to for Label encoding
13
14    print(colored("\n\nColumns Requiring Encoding: \n", 'blue', attrs=['bold']))
15          + colored(categorical_column, 'green', attrs=['bold']))
16
17    # Make Empty Dataframe to decode encoded data Later
18    decode_features = pd.DataFrame()
19
20    ##### Employ Scikit-Learn Label encoding to encode object data #####
21    lab_enc = preprocessing.LabelEncoder()
22    for col in categorical_column:
23        n[col] = lab_enc.fit_transform(n[col])
24        name_mapping = dict(zip(lab_enc.classes_, lab_enc.transform(lab_enc.classes_)))
25
26        ##### Decode Encoded Data #####
27        feature_df = pd.DataFrame([name_mapping])
28        feature_df = feature_df.astype(str)
29        feature_df = (col + "_" + feature_df.iloc[0:])
30        feature_df["Feature"] = col
31        decode_features = decode_features.append(feature_df) # Append Dictionaries to Empty Dataframe for Later Decoding
32
33        ##### Print Encoded Data #####
34        print(colored("Feature: \n", 'blue', attrs=['bold']))
35              + colored(col, 'red', attrs=['bold'])
36              + colored("\nMapping: \n", 'blue', attrs=['bold'])
37              + colored(name_mapping, 'green', attrs=['bold'])
38              + colored("\n\nType n: ", 'blue', attrs=['bold'])
39              + colored(type(n), 'magenta', attrs=['bold'])
40        )
41    n.head(3)
42
43    ##### 2) Make Decoded Factor Dataframe with Description #####
44    factor_list = decode_features.T # Transpose Dataframe and place in new dataframe
45    factor_list = factor_list.replace(np.nan, "/") # nan values with forward slash
46    factor_list["Factors"] = factor_list.astype(str).agg("".join,axis=1).replace(r'[^\\w\\s]|/', '', regex=True) # Aggregate A
47    factor_list.reset_index() # Reset index before copying/assigning it to a new column
48    factor_list["Description"] = factor_list.index # Assign index to column
49
50    return n, factor_list

```

```

1 n, factor_list = encoder(df)
2 print(colored("\n\nEncoded Dataframe\n", 'blue', attrs=['bold']))
3       + colored(type(n), 'blue', attrs=['bold']))
4
5 df = n.copy()

```

Columns Requiring Encoding:

```

['MS Zoning', 'Street', 'Alley', 'Lot Shape', 'Land Contour', 'Utilities', 'Lot Config', 'Land Slope', 'Neighborhood', 'Con
dition 1', 'Condition 2', 'Bldg Type', 'House Style', 'Roof Style', 'Roof Matl', 'Exterior 1st', 'Exterior 2nd', 'Mas Vnr T
ype', 'Exter Qual', 'Exter Cond', 'Foundation', 'Bsmt Qual', 'Bsmt Cond', 'Bsmt Exposure', 'BsmtFin Type 1', 'BsmtFin Type
2', 'Heating', 'Heating QC', 'Central Air', 'Electrical', 'Kitchen Qual', 'Functional', 'Fireplace Qu', 'Garage Type', 'Gar
age Finish', 'Garage Qual', 'Garage Cond', 'Paved Drive', 'Pool QC', 'Fence', 'Misc Feature', 'Sale Type', 'Sale Conditio
n']

```

Feature:

MS Zoning

Mapping:

```
{'A (agr)': 0, 'C (all)': 1, 'FV': 2, 'I (all)': 3, 'RH': 4, 'RL': 5, 'RM': 6}
```

Type n: <class 'pandas.core.frame.DataFrame'>

Feature:

Street

Mapping:

```
{'Grvl': 0, 'Pave': 1}
```

2) Subsequently, multicollinearity was managed by another automated method, as follows:

Method to Eliminate Columns with High Multicollinearity

- 1) Calculate Variance Inflation Factor
- 2) Delete features with VIFs above 2.4 but with no significant relationship with target variable
- 3) Keep features VIFs above 2.4 but with significant relationship with target variable to avoid information loss

Multicollinearity refers to correlation between two or more independent variables which increases standard error (precision of the estimate) of the coefficient. Hence, features exhibiting high multicollinearity can overinflate standard error, thereby, decreasing precision of the estimate. While multicollinearity enlarges model variance, it also expands model dimensions without necessarily enhancing information and so distorts model explainability.

Multicollinearity can easily be computed by the variance inflation factor (VIF) which not only picks out correlation between independent variables but also strength of these correlations. Although, most research papers regard a VIF > 10 as an indicator of strong multicollinearity, nevertheless, there some scholars suggest to select a more cautious threshold of 2.5 which can signal considerable collinearity. Accordingly, Ames House Prediction model implemented a conservative VIF threshold of 2.5 with low correlation to target variable.

```

1 def vifs(*name):
2     n = name # Extract Dataframe by Name...this will create a 3d tuple
3     n = (n[0]) # Convert Tuple to To Dataframe
4
5     analyze, stats = analysis(n) # Call function 'analysis'
6     analyze = analyze[(analyze.Columns != target)] # Remove target Column function 'analysis'
7
8     df_name = [x for x in globals() if globals()[x] is n][0] # Extract Name of Imported Dataframe
9
10    vifs = pd.Series(np.linalg.inv(n.corr().to_numpy()).diagonal(),
11                    index=n.columns,
12                    name='VIF')
13
14    vifs = vifs.drop([target]); # Remove Target 'SalePrice' Column
15    vifs = vifs.to_frame()
16    vifs['Columns'] = vifs.index
17    vifs = vifs.sort_values('VIF', ascending=False)
18
19    # Merge with Analysis to get Correlation with Target Variable
20    vifs = pd.merge(vifs, analyze, on='Columns', how='left')
21    vifs = vifs[(vifs['Corr_Sales'] < 0) & (vifs['VIF']>2.4)]
22
23    vifs = vifs.reset_index()
24    vifs = vifs.sort_values('Columns', ascending=True)
25    drop1 = vifs.Columns.values.tolist()
26
27    return drop1

```

- 3) Using an additional automated method, statistical measures were then employed with supervised filter-based feature selection technique. Firstly, aforementioned “vif” method was called to identify columns exhibiting high multicollinearity but low correlation with sales. Then, columns with less unique features were marked. Consequently, highly skewed columns with low correlation to target were also enlisted. Lastly, Columns with extremely high null values were also keyed out. The lists were then combined to filter these columns out of the data-frame.

Method to Drop Columns

- 1) With High Multicollinearity but Low Correlation to Target
- 2) With Uniques < 2
- 3) With High Skewness and Low Correlation to Target
- 4) Drop Columns With High Nan Values

```

1 def drop_cols(*name):
2     n = name # Extract Dataframe by Name...this will create a 3d tuple
3     n = (n[0]) # Convert Tuple to To Dataframe
4     df_name = [x for x in globals() if globals()[x] is n][0] # Extract Name of Imported Dataframe
5
6     # 1) Drop Columns with With High Multicollinearity & Low Correlation to Sales
7     drop1 = vifs(n) # Call function 'vifs'
8
9     analyze, stats = analysis(n) # Call function 'analysis'
10    print(analyze)
11
12    n = n.fillna(0) #Fill Remaining Missing Values with Zero
13    # Find Mean of Null, Nan and Zero Values Before Any Drops
14    m0 = n.isin([' ', 'NULL', 'NaN', 0]).mean().sort_values(axis=0, ascending=False, inplace=False, kind='quicksort', na_posit
15
16    # 2) Drop Columns with Unique Values Less than threshold
17    unique_counts = pd.DataFrame.from_records([(col, n[col].nunique()) for col in n.columns], # get unique counts
18        columns=['Column_Name', 'Unique']).sort_values(by=['Unique'])
19    unique = unique_counts[unique_counts['Unique'] < 2] #If threshold is Less than 2 then
20    drop2 = [unique['Column_Name'].tolist()] # First List of columns to drop
21
22    # 3) Drop Highly Skewed & Low Sales Correlation Columns OR Low Sales Correlation Columns
23    drop3 = analyze[(analyze['Corr_Sales'] != 1) & (analyze['skewness'] > 0) & (analyze['Corr_Sales'] < 0) | (analyze['Corr_
24
25    print(colored("\nDrop 3: \n ", 'blue', attrs=['bold']))
26    +colored(drop3, 'magenta', attrs=['bold']))
27
28    drop3 = drop3.sort_values(by='Columns', ascending=True)
29    drop3 = drop3['Columns'].tolist() # Second List of columns to drop
30
31    drop = drop1 + drop2 + drop3 + delete_features # Final List of columns to drop
32    print(drop)
33    n = n.drop(drop,1) # Drop Columns using Final List of columns to drop
34
35    # Find Mean of Null, Nan and Zero Values Before Dropping
36    m1 = n.isin([' ', 'NULL', 'NaN', 0]).mean().sort_values(axis=0, ascending=False, inplace=False, kind='quicksort', na_posit
37
38    # 4) Drop Columns With High Nan Values
39    drop_thresh = .90 # Identify Drop Threshold
40    n = n.loc[:, df.isin([' ', 'NULL', 'NaN', 0]).mean() < drop_thresh] # drop columns if Mean is > 0.90
41
42    #df = df.fillna(0) #Fill Remaining Missing Values with Zero
43    n = n.replace(["NaN"], 0).sort_values(by=target, ascending=False) # Replace all Nan Values with Zero
44
45    # Find Mean of Null, Nan and Zero Values After Dropping
46    m2 = n.isin([' ', 'NULL', 'NaN']).mean().sort_values(axis=0, ascending=False, inplace=False, kind='quicksort', na_positior
47
48    #Print Results
49    print(colored("\nDataframe Average Null Values Before Any Drops\n", 'blue', attrs=['bold']))
50    +colored(m0, 'magenta', attrs=['bold'])
51    +colored("\n\n Low Correlation Columns to Drop: ", 'green', attrs=['bold'])
52    + colored(drop1, 'red', attrs=['bold'])
53    +colored("\n\nDataframe Average Null Values After Low Correlation Columns Drop\n", 'green', attrs=['bold'])
54    +colored(m1, 'red', attrs=['bold'])
55    +colored("\n\n Drop Columns if Mean is > 0.90 \n", 'green', attrs=['bold'])
56    + colored("\nDataframe Average Null Values After Drop and 'Nan' Replacement\n", 'blue', attrs=['bold'])
57    +colored(m2, 'magenta', attrs=['bold'])
58    +colored(type(m2), 'magenta', attrs=['bold'])
59    )
60    return n
61
62    # Return Function
63    n = drop_cols(df)
64    df = n.copy()

```

- 4) Lastly, an automated method was employed to replace outliers with mode, that is, most frequent value.

Method to Explore and Adjust Outliers

Replace Outlier Values with Mode (Most Frequent Value)

```

1 def outliers(*name):
2     n = name # Extract Dataframe by Name...this will create a 3d tuple
3     n = (n[0]) # Convert Tuple to DataFrame
4     df_name = [x for x in globals() if globals()[x] is n][0] # Extract Name of Imported Dataframe
5
6     cols = n.columns # ALL Columns
7
8     # Numeric Columns
9     numerics = ['int16', 'int32', 'int64', 'float16', 'float32', 'float64']
10    numeric_cols = n.select_dtypes(include=numerics)
11    numeric_cols = numeric_cols.columns.tolist()
12
13    # Object Columns
14    categorical_cols = list(set(cols) - set(numeric_cols))
15
16    # Skewed Columns
17    skewed_cols = analyze[(analyze['skewness'] > 0) | (analyze['skewness'] < 0)]
18    skewed_cols = skewed_cols['Columns'].tolist()
19
20    # Replace Outliers
21    for col in n.columns:
22        if col in skewed_cols:
23            print(colored(col, 'magenta', attrs=['bold']))
24            + colored(" is Skewed... ", 'blue', attrs=['bold'])
25            )
26
27            mode = n[col].mode()
28            mode = mode[0]
29
30            if col in numeric_cols:
31                print(colored(col, 'magenta', attrs=['bold']))
32                + colored(" Column Type is: ", 'blue', attrs=['bold'])
33                + colored(n[col].dtypes, 'red', attrs=['bold'])
34                )
35                #Calculate quantiles and IQR
36                Q1 = n[col].quantile(0.25) # Same as np.percentile but maps (0,1) and not (0,100)
37                Q3 = n[col].quantile(0.75)
38                IQR = Q3 - Q1
39                # Replace with Mode
40                n[col] = np.where((n[col] < (Q1 - 1.5 * IQR)) | (n[col] > (Q3 + 1.5 * IQR)), mode, n[col])
41
42                print(colored("\nReplaced ", 'blue', attrs=['bold']))
43                + colored(col, 'magenta', attrs=['bold'])
44                + colored(" Skewed Values by Mode: ", 'blue', attrs=['bold'])
45                + colored(mode, 'red', attrs=['bold'])
46                + colored("\n", 'magenta', attrs=['bold'])
47                + colored((n[col]), 'green', attrs=['bold'])
48                )
49
50            else:
51                print("")
52    df_transformed = n.copy()
53    return df_transformed

```

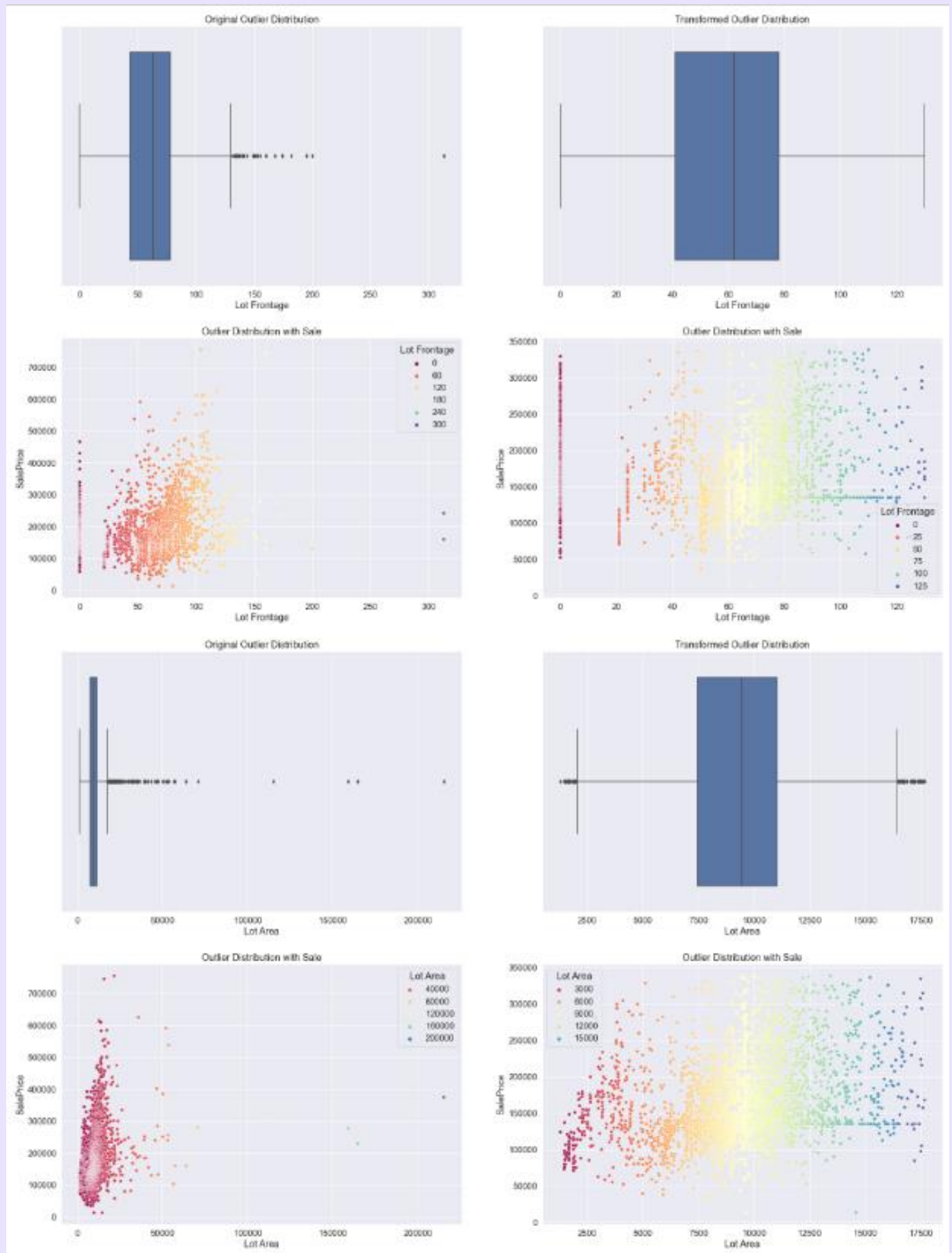

This was followed by another method to present graphical illustration of original and adjusted data side by side:

Visualize Data Distribution of both Original and New Dataframe

```

1 df_numeric = df_transformed.select_dtypes(exclude='object')
2 for col in df_numeric.columns: # Iterate over each Column and Create Visuals
3     ##### New Figure #####
4     figa = plt.figure(figsize=(30, 20))
5     sns.set(font_scale=1.5)
6
7     fig1 = figa.add_subplot(221); sns.boxplot(df[col])
8     fig1 = plt.title('Original Outlier Distribution')
9
10    fig2 = figa.add_subplot(222); sns.boxplot(df_transformed[col])
11    fig2 = plt.title('Transformed Outlier Distribution')
12
13    fig3 = figa.add_subplot(223);
14    sns.scatterplot(x = df[col], y = df[target], hue=df[col], palette= 'Spectral')
15    fig3 = plt.title('Outlier Distribution with Sale')
16
17    fig4 = figa.add_subplot(224);
18    sns.scatterplot(x = df_transformed[col], y = df_transformed[target], hue=df_transformed[col], palette= 'Spectral')
19    fig4 = plt.title('Outlier Distribution with Sale')
20
21    ##### New Figure #####
22    figb = plt.figure(figsize=(20, 10))
23    sns.set(font_scale=1.5)
24
25    fig5 = figb.add_subplot(221);
26    fig5 = sns.distplot(df[target][~df[target].isnull()], axlabel="Nor. Dist.", fit=st.norm, fit_kws={"color":"red"})
27    fig5 = plt.title('Distribution of Sales Price')
28    (mu5, sigma5) = st.norm.fit(df[target])
29    fig5 = plt.legend(['Normal Distribution \n ($\mu=$ {:.2f} and $\sigma=$ {:.2f} )'.format(mu5, sigma5)], loc='best', fancybo
30
31    fig6 = figb.add_subplot(222);
32    fig6 = sns.distplot(df_transformed[target][~df_transformed[target].isnull()], axlabel="Nor. Dist.", fit=st.norm, fit_kws={"c
33    fig6 = plt.title('Distribution of Sales Price')
34    (mu6, sigma6) = st.norm.fit(df_transformed[target])
35    fig6 = plt.legend(['Normal Distribution \n ($\mu=$ {:.2f} and $\sigma=$ {:.2f} )'.format(mu6, sigma6)], loc='best', fancybo
36
37    ##### New Figure #####
38    figc = plt.figure(figsize=(20, 10))
39    sns.set(font_scale=1.5)
40    fig7 = figc.add_subplot(221);
41    fig7 = st.probplot(df[target][~df[target].isnull()], plot=plt)
42    fig7 = plt.title('Probability Plot')
43
44    fig8 = figc.add_subplot(222);
45    fig8 = st.probplot(df_transformed[target][~df_transformed[target].isnull()], plot=plt)
46    fig8 = plt.title('Probability Plot')
47

```



5) Summary of Training Different Regression Models

5a) Machine Learning Regression Algorithm Development Approach

Using random search, an automated optimal hyper-parameter search method was created to find optimal model parameters. This approach was employed because best hyperparameters are not automatically learnt within estimators and its manual search not only slows down model development but may also lead to ineffective model construction. Hence, an exhaustive random search approach was used to pass parameter arguments to the constructor in order to find optimal hyperparameters for each model, as shown below:

RIDGE REGRESSION MODELS

Random Search Method to Find 'Best Parameters'to 'Build Ridge Regression Model WITH Optimal Hyperparameters'

```

1 def random_search_r(X_train, y_train):
2
3     # Define Evaluation
4     cv = RepeatedKFold(n_splits=10, n_repeats=3, random_state=1)
5
6     # Define Search Space
7     space = dict()
8     space['solver'] = ['auto', 'svd', 'cholesky', 'lsqr', 'sparse_cg', 'sag', 'saga', 'lbfgs']
9     space['alpha'] = loguniform(1e-5, 100)
10    space['fit_intercept'] = [True, False]
11    space['normalize'] = [True, False]
12    space['max_iter'] = [500, 1000, 1500]
13
14    # Define Model
15    ridge_model = Ridge(random_state=rs, max_iter=1000)
16
17    # Define Search
18    search = RandomizedSearchCV(ridge_model, space, n_iter=500, scoring='neg_mean_absolute_error', n_jobs=-1, cv=cv, random_
19
20    # Execute Search
21    result_r = search.fit(X_train, y_train)
22
23    # Summarize Result
24    best_score_r = result_r.best_score_
25
26    best_params_r = result_r.best_params_
27    best_params_r["best_score"] = best_score_r # Add 'best_score' to 'best_params' Dictionary
28
29    best_params_r = pd.DataFrame([best_params_r]) # Dictionary To dataframe
30
31    # Get Optimal Variables
32    opt_alpha_r = best_params_r['alpha'].iloc[0]
33    opt_alpha_r = '{:.6f}'.format(opt_alpha_r)
34    opt_alpha_r = float(opt_alpha_r) # Back to Float
35
36    opt_solver_r = best_params_r['solver'].iloc[0]
37    opt_fit_intercept_r = best_params_r['fit_intercept'].iloc[0]
38    opt_normalize_r = best_params_r['normalize'].iloc[0]
39    opt_max_iter_r = best_params_r['max_iter'].iloc[0]
40
41    # Define Optimal Parameters
42    optimal_params_r = {'solver': opt_solver_r,
43                        'alpha': opt_alpha_r,
44                        'fit_intercept': opt_fit_intercept_r,
45                        'normalize': opt_normalize_r,
46                        'max_iter': opt_max_iter_r,
47                        'random_state': rs}
48
49    return best_params_r, optimal_params_r
50
51 best_params_r, optimal_params_r = random_search_r(X_train, y_train) # Call Method 'random_search_r' to get optimal hyperpara
52 best_params_r

```

	alpha	fit_intercept	max_iter	normalize	solver	best_score
0	75.322052	True	500	False	auto	-23514.379109

LAGSSO REGRESSION MODELS

Random Search Method to Find 'Best Parameters'to 'Build Lasso Regression Model WITH Optimal Hyperparameters'

```

1 def random_search_1(X_train, y_train):
2
3     # Define Evaluation
4     cv = RepeatedKfold(n_splits=10, n_repeats=3, random_state=1)
5
6     # Define Search Space
7     space = dict()
8     space['alpha'] = loguniform(1e-5, 100)
9     space['fit_intercept'] = [True, False]
10    space['normalize'] = [True, False]
11    space['precompute'] = [True, False]
12    space['tol'] = loguniform(1e-4, 100)
13    space['selection'] = ['cyclic', 'random']
14
15    # Define Model
16    lasso_model = Lasso(random_state=rs, max_iter=1000)
17
18    # Define Search
19    search = RandomizedSearchCV(lasso_model, space, n_iter=1000, scoring='neg_mean_absolute_error', n_jobs=-1, cv=cv, random
20
21    # Execute Search
22    result_1 = search.fit(X_train, y_train)
23
24    # Summarize Result
25    best_score_1 = result_1.best_score_
26    best_params_1 = result_1.best_params_
27    best_params_1["best_score"] = best_score_1 # Add 'best_score' to 'best_params' Dictionary
28
29    best_params_1 = pd.DataFrame([best_params_1]) # Dictionary To dataframe
30
31    # Get Optimal Variables
32    opt_alpha_1 = best_params_1['alpha'].iloc[0]
33    opt_fit_intercept_1 = best_params_1['fit_intercept'].iloc[0]
34    opt_normalize_1 = best_params_1['normalize'].iloc[0]
35    opt_precompute_1 = best_params_1['precompute'].iloc[0]
36    opt_tol_1 = best_params_1['tol'].iloc[0]
37    opt_selection_1 = best_params_1['selection'].iloc[0]
38
39    # Define Optimal Parameters
40    optimal_params_1 = {'alpha': opt_alpha_1,
41                        'fit_intercept': opt_fit_intercept_1,
42                        'normalize': opt_normalize_1,
43                        'precompute': opt_precompute_1,
44                        'tol': opt_tol_1,
45                        'selection': opt_selection_1}
46
47
48    return best_params_1, optimal_params_1
49
50 best_params_1, optimal_params_1 = random_search_1(X_train, y_train) # Call Method 'random_search_1' to get optimal hyperpara
51 best_params_1

```

	alpha	fit_intercept	normalize	precompute	selection	tol	best_score
0	9.126019	True	True	False	random	0.031602	{'alpha': 9.126018989876911, 'fit_intercept': ...

ELASTIC-NET REGRESSION MODELS

Randomized Search Method to Find 'Best Parameters' to 'Build Elastic Net Regression Model WITH Optimal Hyperparameters'

```

1 def random_search_en(X_train, y_train):
2
3     # Define Evaluation
4     cv = RepeatedKfold(n_splits=10, n_repeats=3, random_state=1)
5
6     # Define Search Space
7     space = dict()
8     space['alpha'] = loguniform(1.0, 1.5, 3, 3.5)
9     space['l1_ratio'] = loguniform(0.5, 1)
10    space['fit_intercept'] = [True, False]
11    space['normalize'] = [True, False]
12    space['precompute'] = [True, False]
13    space['copy_X'] = [True, False]
14    space['tol'] = [1e-4, 1e-6, 1e-9]
15    space['warm_start'] = [True, False]
16    space['positive'] = [True, False]
17    space['selection'] = ['cyclic', 'random']
18
19    # Define Model
20    en_model = ElasticNet(random_state=rs, max_iter=1000)
21
22    # Define Search
23    search_en = RandomizedSearchCV(en_model, space, n_iter=1000, scoring='neg_mean_absolute_error', n_jobs=-1, cv=cv, random
24
25    # Execute Search
26    result_en = search_en.fit(X_train, y_train)
27
28    # Summarize result
29    best_score_en = result_en.best_score_
30
31    best_params_en = result_en.best_params_
32    best_params_en["best_score"] = best_score_en # Add 'best_score' to 'best_params' Dictionary
33    best_params_en = pd.DataFrame([best_params_en]) # Dictionary To dataframe
34
35    # Get Optimal Variables
36    opt_alpha_en = best_params_en['alpha'].iloc[0]
37    opt_l1_ratio_en = best_params_en['l1_ratio'].iloc[0]
38    opt_fit_intercept_en = best_params_en['fit_intercept'].iloc[0]
39    opt_normalize_en = best_params_en['normalize'].iloc[0]
40    opt_precompute_en = best_params_en['precompute'].iloc[0]
41    opt_copy_X_en = best_params_en['copy_X'].iloc[0]
42    opt_tol_en = best_params_en['tol'].iloc[0]
43    opt_tol_en = best_params_en['tol'].iloc[0]
44    opt_warm_start_en = best_params_en['warm_start'].iloc[0]
45    opt_positive_en = best_params_en['positive'].iloc[0]
46    opt_selection_en = best_params_en['selection'].iloc[0]
47
48    # Define Optimal Parameters
49    optimal_params_en = {'alpha': opt_alpha_en,
50                        'l1_ratio': opt_l1_ratio_en,
51                        'fit_intercept': opt_fit_intercept_en,
52                        'normalize': opt_normalize_en,
53                        'precompute': opt_precompute_en,
54                        'copy_X': opt_copy_X_en,
55                        'tol': opt_tol_en,
56                        'warm_start': opt_warm_start_en,
57                        'positive': opt_positive_en,
58                        'selection': opt_selection_en}
59    return best_params_en, optimal_params_en
60
61 best_params_en, optimal_params_en = random_search_en(X_train, y_train) # Call Function 'random_search_en' to get optimal hyp
62 best_params_en

```

	alpha	copy_X	fit_intercept	l1_ratio	normalize	positive	precompute	selection	tol	warm_start	best_score
0	6.9506	True	True	0.967434	False	True	False	random	1.000000e-09	True	{'alpha': 6.950599680152801, 'copy_X': True, '...

XGBOOST REGRESSION MODELS ¶

Random Search Method to Find 'Best Parameters'to 'Build XGBoost Regression Model WITH Optimal Hyperparameters'

```

1 def random_search_xgb(X_train, y_train):
2
3     # Define Evaluation
4     cv = RepeatedKFold(n_splits=10, n_repeats=3, random_state=1)
5
6     # Define Search Space
7     space = dict()
8     space['learning_rate'] = [0.1, 20]
9     space['subsample'] = [0.0, 1.0]
10    space['criterion'] = ['friedman_mse', 'squared_error', 'mse']
11    space['max_features'] = ['auto', 'sqrt', 'log2']
12
13    # Define Model
14    xgb_model = GradientBoostingRegressor(random_state=rs)
15
16    # Define Search
17    search_xgb = RandomizedSearchCV(xgb_model, space, n_iter=1000, scoring='neg_mean_absolute_error', n_jobs=-1, cv=cv, rand
18
19    # Execute Search
20    result_xgb = search_xgb.fit(X_train, y_train)
21
22    # Summarize Result
23    best_score_xgb = result_xgb.best_score_
24    best_params_xgb = result_xgb.best_params_
25    best_params_xgb["best_score"] = best_score_xgb # Add 'best_score' to 'best_params' Dictionary
26
27    best_params_xgb = pd.DataFrame([best_params_xgb]) # Dictionary To dataframe
28
29    # Get Optimal Variables
30    opt_learning_rate_xgb = best_params_xgb['learning_rate'].iloc[0]
31    opt_subsample_xgb = best_params_xgb['subsample'].iloc[0]
32    opt_criterion_xgb = best_params_xgb['criterion'].iloc[0]
33    opt_max_features_xgb = best_params_xgb['max_features'].iloc[0]
34
35    # Define Optimal Parameters
36    optimal_params_xgb = {'learning_rate': opt_learning_rate_xgb,
37                          'subsample': opt_subsample_xgb,
38                          'criterion': opt_criterion_xgb,
39                          'max_features': opt_max_features_xgb}
40
41    return best_params_xgb, optimal_params_xgb
42
43 best_params_xgb, optimal_params_xgb = random_search_xgb(X_train, y_train) # Call Function 'random_search_xgb' to get optimal
44 best_params_xgb

```

	subsample	max_features	learning_rate	criterion	best_score
0	1.0	auto	0.1	friedman_mse	{'subsample': 1.0, 'max_features': 'auto', 'le...

5b) Summarizing Employed Models

Following four main regression models have been used to predict house prices.

1) Ridge Regression (RR) Models

Due to the dependant nature of multiple variables in predicting variable 'y' where the output is influenced by multiple factors, ridge regression algorithm was employed. A single method was created to measure predictive capability of the following two RR models:

Method to 'Build Ridge Regression WITH & Without Optimal Hyperparameters'

```

1  # Build a Regression model with Optimal Class Weights
2  def build_op_r(X_train, y_train, X_test, threshold=0.5, best_params=None):
3
4      model = Ridge(random_state=rs, max_iter = 1000)
5
6      # If best parameters are provided
7      if best_params:
8          model = make_pipeline(RobustScaler(),
9                                Ridge(solver = best_params_r['solver'].iloc[0],
10                                     alpha = best_params_r['alpha'].iloc[0],
11                                     fit_intercept = best_params_r['fit_intercept'].iloc[0],
12                                     normalize = best_params_r['normalize'].iloc[0],
13                                     max_iter = best_params_r['max_iter'].iloc[0],
14                                     random_state=rs
15                                )
16                                )
17      model.fit(X_train, y_train)
18
19      # Get Prediction
20      pred = model.predict(X_test)
21      #pred = pred * (y_train.std()) + y_train.mean()
22      test_ids = X_test.index
23      pred = pd.DataFrame({'Id': test_ids, 'Predicted Price': pred})
24
25      # Get Model Directory
26      print(colored("Model Directory: \n", 'cyan', attrs=['bold'])
27            + colored(dir(model), 'magenta', attrs=['bold']))
28
29      # Get Model Features
30      feature_names = model[-1].get_feature_names_out()
31      X_train_preprocessed = pd.DataFrame(model[-1].transform(X_train), columns=feature_names) # X_train to compute stand
32      features = pd.DataFrame(model[-1].coef_ * X_train_preprocessed.std(axis=0), columns=["Coefficients"], index=feature_
33
34
35  else:
36      model.fit(X_train, y_train)
37
38      # Get Prediction
39      pred = model.predict(X_test)
40      #pred = pred * (y_train.std()) + y_train.mean()
41      test_ids = X_test.index
42      pred = pd.DataFrame({'Id': test_ids, 'Predicted Price': pred})
43
44      # Get Model Directory
45      print(colored("Model Directory: \n", 'cyan', attrs=['bold'])
46            + colored(dir(model), 'magenta', attrs=['bold']))
47
48      # Get Model Features
49      feature_names = model.check_feature_names
50      X_train_preprocessed = X_train.copy() #.T # Transpose X_train to compute standard deviation of the related feature
51
52      # Multiply model coefficients by std of features to reduce all coefficients to same unit of measure.
53      features = pd.DataFrame(model.coef_ * X_train_preprocessed.std(axis=0), columns=["Coefficients"], index=X_train.co
54
55  features['Features'] = features.index
56  #features = features.loc[features['Coefficients'] > 0] # Uncomment if only features with Coeff > 0 are required
57  features = features.sort_values(by=['Coefficients'], ascending=False).reset_index(drop=True)
58
59  # Plot Features
60  data = features.copy()
61  data = data[data["Coefficients"] != 0]
62  data.set_index('Features', inplace=True)
63  data.plot.barh(figsize=(30,10), color='green')
64
65  title = "Feature Importance in Predicting " + target
66  plt.title(title)
67  plt.axvline(x=0, color="red")
68  plt.xlabel("Coefficient Values Corrected by Features' std. dev.")
69  plt.subplots_adjust(left=0.3)
70
71  # Get Model Results
72  res = result(model)
73
74  return model, res, features, pred

```

1a) RR Model 1 Without Optimal Hyperparameter Tuning: A simple algorithm was created without optimal hyperparameter tuning.

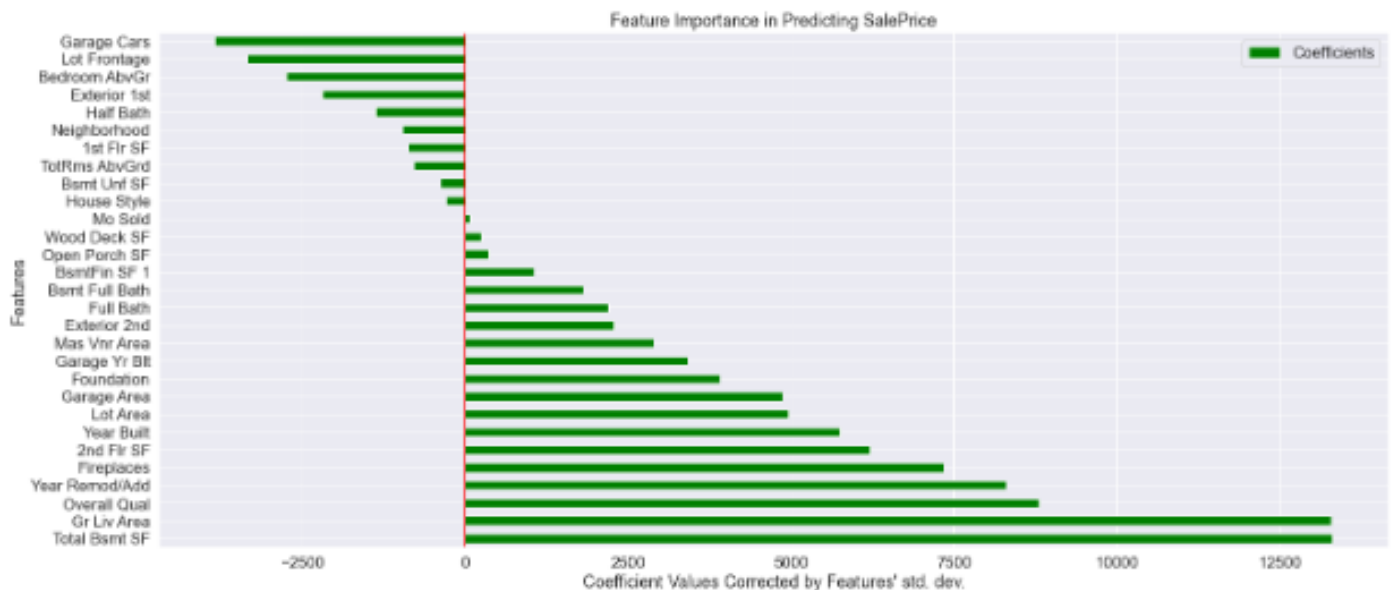
Model 1: RIDGE Without 'Optimal Hypertuned Parameters'

```
1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 model1, res1, features1, pred1 = build_op_r(X_train, y_train, X_test, best_params=None) # Call Method to 'Build Logistic Reg
3
4
5 #Collect & Append Results
6 r1 = np.array2string(res1)####
7 r1 = pd.DataFrame([r1.split(';') for x in r1.split('\n')])
8 r1 = r1.rename(columns={0: 'SCORE'})
9 r1['MODEL'] = 'Regression'
10 r1['MODEL#'] = 'Model 1'
11 r1
```

Model Directory:

```
[ '_abstractmethods_', '_class_', '_delattr_', '_dict_', '_dir_', '_doc_', '_eq_', '_format_', '_ge_', '_geta
ttribute_', '_getstate_', '_gt_', '_hash_', '_init_', '_init_subclass_', '_le_', '_lt_', '_module_', '_ne_',
'_new_', '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_setstate_', '_sizeof_', '_str_', '_subclasshook_',
'_weakref_', '_abc_impl_', '_check_feature_names_', '_check_n_features_', '_decision_function_', '_estimator_type_', '_get_par
am_names_', '_get_tags_', '_more_tags_', '_normalize_', '_repr_html_', '_repr_html_inner_', '_repr_mimebundle_', '_set_intercept_',
'_validate_data_', '_alpha_', '_coef_', '_copy_X_', '_feature_names_in_', '_fit_', '_fit_intercept_', '_get_params_', '_intercept_', '_max_ite
r_', '_n_features_in_', '_n_iter_', '_normalize_', '_positive_', '_predict_', '_random_state_', '_score_', '_set_params_', '_solver_', '_tol_' ]
```

	SCORE	MODEL	MODEL#
0	-1.30784598e+09	Regression	Model 1



1b) RR Model 2 WITH Optimal Hyperparameter Tuning: A modified algorithm with auto-hyper-tuned parameters was created using random search method described above.

Model 2: RIDGE With 'Optimal Hypertuned Parameters'

```

1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 model2, res2, features2, pred2 = build_op_r(X_train, y_train, X_test, best_params=optimal_params_r) # Call Method to 'Build
3
4 #Collect & Append Results
5 r2 = np.array2string(res2)####
6 r2 = pd.DataFrame([r2.split(';') for x in r2.split('\n')])
7 r2 = r2.rename(columns={0: 'SCORE'})
8 r2['MODEL'] = 'Regression Optimal'
9 r2['MODEL#'] = 'Model 2'
10 r2

```

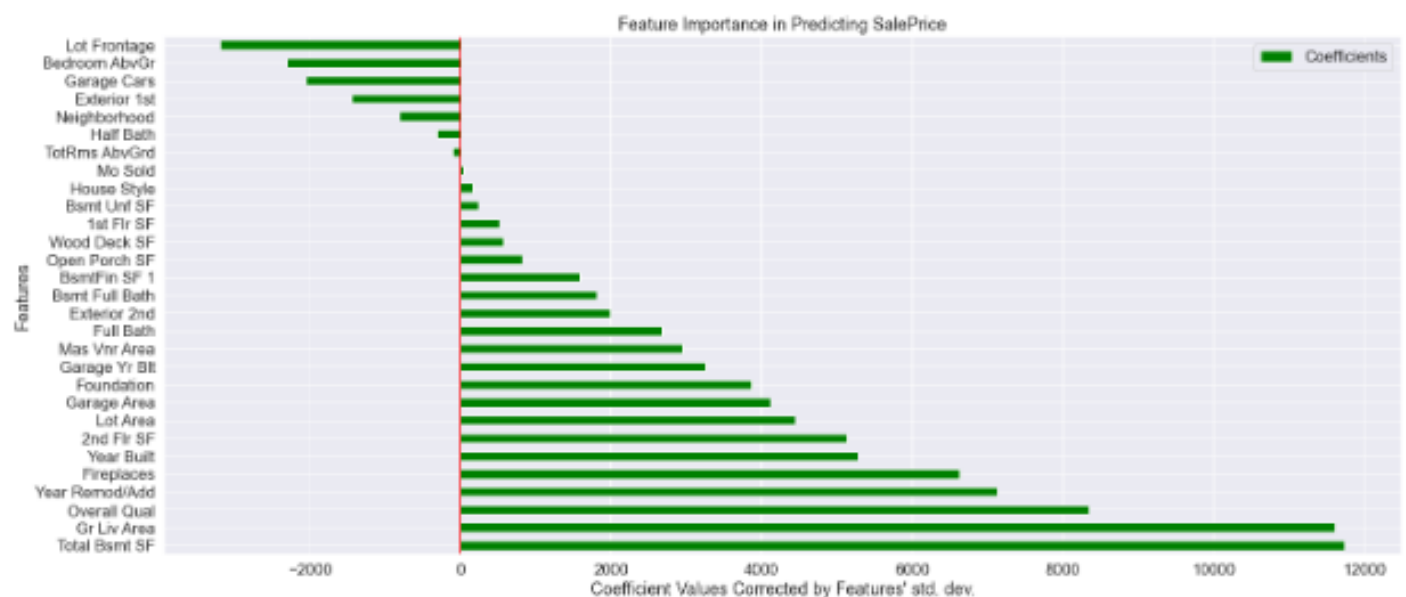
Model Directory:

```

['_abstractmethod_', '_annotations_', '_class_', '_delattr_', '_dict_', '_dir_', '_doc_', '_eq_', '_format_',
'_ge_', '_getattr_', '_getitem_', '_getstate_', '_gt_', '_hash_', '_init_', '_init_subclass_', '_le_',
'_len_', '_lt_', '_module_', '_ne_', '_new_', '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_setsta',
'te_', '_sizeof_', '_sklearn_is_fitted_', '_str_', '_subclasshook_', '_weakref_', '_abc_impl_', '_can_inverse_transfor',
'm_', '_can_transform_', '_check_feature_names_', '_check_fit_params_', '_check_n_features_', '_estimator_type_', '_final_estimator_',
'_fit_', '_get_param_names_', '_get_params_', '_get_tags_', '_iter_', '_log_message_', '_more_tags_', '_replace_estimator_', '_repr_htm',
'l_', '_repr_html_inner_', '_repr_mimebundle_', '_required_parameters_', '_set_params_', '_sk_visual_block_', '_validate_data_', '_v',
'alidate_names_', '_validate_steps_', 'classes_', 'decision_function_', 'feature_names_in_', 'fit_', 'fit_predict_', 'fit_transform_',
'get_feature_names_out_', 'get_params_', 'inverse_transform_', 'memory_', 'n_features_in_', 'named_steps_', 'predict_', 'predict_log',
'proba_', 'predict_proba_', 'score_', 'score_samples_', 'set_params_', 'steps_', 'transform_', 'verbose_']

```

	SCORE	MODEL	MODEL#
0	-1.30393778e+09	Regression Optimal	Model 2



2) Lasso Regression (LR) Models

A single method was employed to run lasso regression models with and without optimal hyperparameter tuning:

Method to 'Build LASSO Regression WITH & Without Optimal Hyperparameters'

```

1  # Build a LASSO model with Optimal Class Weights
2  def build_op_l(X_train, y_train, X_test, threshold=0.5, best_params=None):
3
4      model = Ridge(random_state=rs, max_iter = 1000)
5
6      # If best parameters are provided
7      if best_params:
8          model = make_pipeline(RobustScaler(),
9                                Lasso(alpha = best_params_l['alpha'].iloc[0],
10                                     fit_intercept = best_params_l['fit_intercept'].iloc[0],
11                                     normalize = best_params_l['normalize'].iloc[0],
12                                     precompute = best_params_l['precompute'].iloc[0],
13                                     tol = best_params_l['tol'].iloc[0],
14                                     selection = best_params_l['selection'].iloc[0],
15                                     random_state=rs
16                                )
17
18      model.fit(X_train, y_train)
19
20      # Get Prediction
21      pred = model.predict(X_test)
22      #pred = pred * (y_train.std() + y_train.mean())
23      test_ids = X_test.index
24      pred = pd.DataFrame({'Id': test_ids, 'Predicted Price': pred})
25
26      # Get Model Directory
27      print(colored("Model Directory: \n", 'cyan', attrs=['bold'])
28            + colored(dir(model), 'magenta', attrs=['bold']))
29
30      # Get Model Features
31      feature_names = model[-1].get_feature_names_out() #https://scikit-learn.org/stable/auto_examples/inspection/plot_Li
32      X_train_preprocessed = pd.DataFrame(model[-1].transform(X_train), columns=feature_names) # X_train to compute stand
33      features = pd.DataFrame(model[-1].coef_ * X_train_preprocessed.std(axis=0), columns=["Coefficients"], index=feature
34
35  else:
36      model.fit(X_train, y_train)
37
38      # Get Prediction
39      pred = model.predict(X_test)
40      #pred = pred * (y_train.std() + y_train.mean())
41      test_ids = X_test.index
42      pred = pd.DataFrame({'Id': test_ids, 'Predicted Price': pred})
43
44      # Get Model Directory
45      print(colored("Model Directory: \n", 'cyan', attrs=['bold'])
46            + colored(dir(model), 'magenta', attrs=['bold']))
47
48      # Get Model Features
49      feature_names = model.check_feature_names
50      X_train_preprocessed = X_train.copy() #.T # Transpose X_train to compute standard deviation of the related feature
51      # Multiply model coefficients by std of features to reduce all coefficients to same unit of measure.
52      features = pd.DataFrame(model.coef_ * X_train_preprocessed.std(axis=0), columns=['Coefficients'], index=X_train.col
53
54      features['Features'] = features.index
55      #features = features.loc[features['Coefficients'] > 0] # Uncomment if only features with Coeff > 0 are required
56      features = features.sort_values(by=['Coefficients'], ascending=False).reset_index(drop=True)
57
58      # Plot Features
59      data = features.copy()
60      data = data[data["Coefficients"] != 0]
61      data.set_index('Features', inplace=True)
62      data.plot.barh(figsize=(30,10), color='green')
63
64      title = "Feature Importance in Predicting " + target
65      plt.title(title)
66      plt.axvline(x=0, color="red")
67      plt.xlabel("Coefficient Values Corrected by Features' std. dev.")
68      plt.subplots_adjust(left=0.3)
69
70      # Get Model Results
71      res = result(model)
72
73      return model, res, features, pred

```


2a) LR Model 1 Without Optimal Hyperparameter Tuning: A simple algorithm was created without any hyperparameter tuning.

Model 3: LASSO Without 'Optimal Hypertuned Parameters'

```

1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 model3, res3, features3, pred3 = build_op_l(X_train, y_train, X_test, best_params=None) # Call Method to 'Build Logistic Reg
3
4 #Collect & Append Results
5 r3 = np.array2string(res3)####
6 r3 = pd.DataFrame([r3.split(';') for x in r3.split('\n')])
7 r3 = r3.rename(columns={0: 'SCORE'})
8 r3['MODEL'] = 'Lasso'
9 r3['MODEL#'] = 'Model 3'
10 r3

```

Model Directory:

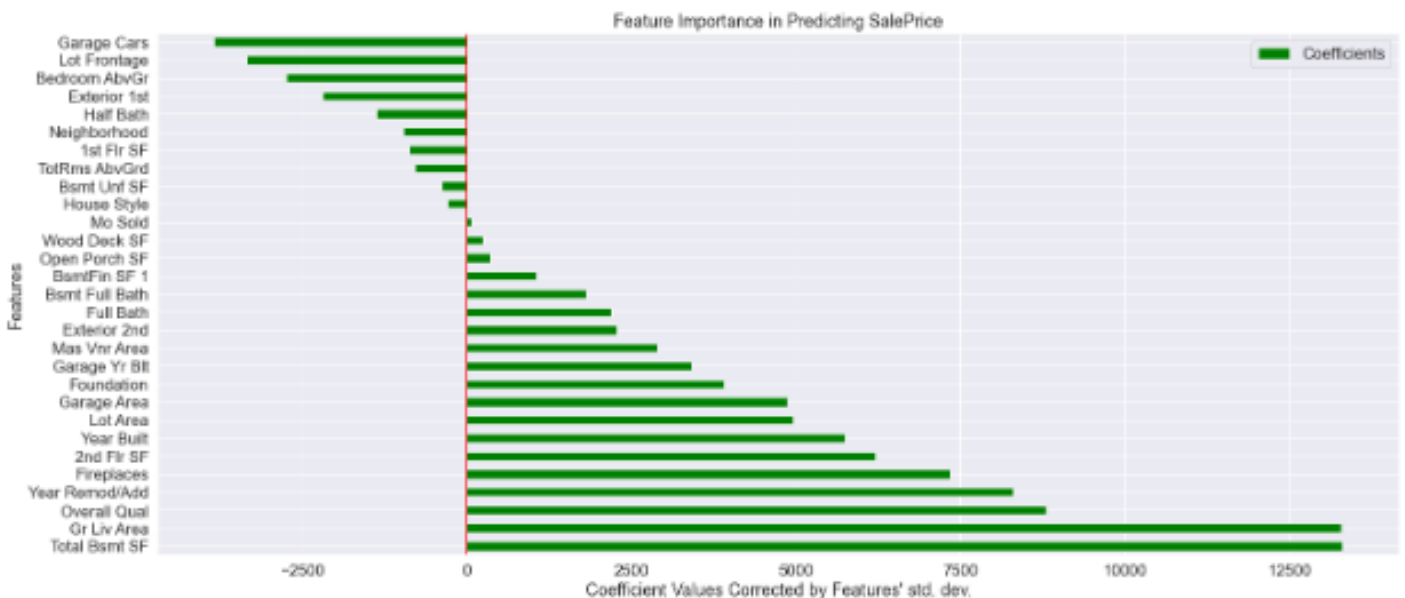
```

['_abstractmethod_', '_class_', '_delattr_', '_dict_', '_dir_', '_doc_', '_eq_', '_format_', '_ge_', '_geta
ttribute_', '_getstate_', '_gt_', '_hash_', '_init_', '_init_subclass_', '_le_', '_lt_', '_module_', '_ne_',
'_new_', '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_setstate_', '_sizeof_', '_str_', '_subclasshook_',
'_weakref_', '_abc_impl', '_check_feature_names', '_check_n_features', '_decision_function', '_estimator_type', '_get_par
am_names', '_get_tags', '_more_tags', '_normalize', '_repr_html_', '_repr_html_inner', '_repr_mimebundle_', '_set_intercept',
'_validate_data', '_alpha', '_coef_', '_copy_X', '_feature_names_in_', '_fit', '_fit_intercept', '_get_params', '_intercept_', '_max_ite
r', '_n_features_in_', '_n_iter_', '_normalize', '_positive', '_predict', '_random_state', '_score', '_set_params', '_solver', '_tol']

```

SCORE MODEL MODEL#

0 -1.30784598e+09 Lasso Model 3



2b) LR Model 2 WITH Optimal Hyperparameter Tuning: A modified algorithm with auto hyper-tuned parameters was created using random search method described above.

Model 4: LASSO With 'Optimal Hypertuned Parameters'

```
1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 model4, res4, features4, pred4 = build_op_l(X_train, y_train, X_test, best_params=optimal_params_1) # Call Method to 'Build
3
4 #Collect & Append Results
5 r4 = np.array2string(res4)####
6 r4 = pd.DataFrame([r4.split(';') for x in r4.split('\n')])
7 r4 = r4.rename(columns={0: 'SCORE'})
8 r4['MODEL'] = 'Lasso Optimal'
9 r4['MODEL#'] = 'Model 4'
10 r4
```

Model Directory:

```
['_abstractmethods_', '_annotations_', '_class_', '_delattr_', '_dict_', '_dir_', '_doc_', '_eq_', '_format_',
'_ge_', '_getattr_', '_getitem_', '_getstate_', '_gt_', '_hash_', '_init_', '_init_subclass_', '_le_',
'_len_', '_lt_', '_module_', '_ne_', '_new_', '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_setsta',
te_', '_sizeof_', '_sklearn_is_fitted_', '_str_', '_subclasshook_', '_weakref_', '_abc_impl_', '_can_inverse_transfor',
m_', '_can_transform_', '_check_feature_names_', '_check_fit_params_', '_check_n_features_', '_estimator_type_', '_final_estimator_',
'_fit_', '_get_param_names_', '_get_params_', '_get_tags_', '_iter_', '_log_message_', '_more_tags_', '_replace_estimator_', '_repr_html',
l_', '_repr_html_inner_', '_repr_mimebundle_', '_required_parameters_', '_set_params_', '_sk_visual_block_', '_validate_data_', '_v',
alidate_names_', '_validate_steps_', '_classes_', '_decision_function_', '_feature_names_in_', '_fit_', '_fit_predict_', '_fit_transform_',
'_get_feature_names_out_', '_get_params_', '_inverse_transform_', '_memory_', '_n_features_in_', '_named_steps_', '_predict_', '_predict_log',
proba_', '_predict_proba_', '_score_', '_score_samples_', '_set_params_', '_steps_', '_transform_', '_verbose_']
```

SCORE	MODEL	MODEL#
0 -1.3033701e+09	Lasso Optimal	Model 4



3) Elastic-Net (EN) Models

A single method was employed to run elastic-net models with and without optimal hyperparameter tuning:

Method to 'Build ELASTIC-NET Regression WITH & Without Optimal Hyperparameters'

```

1  # Build a LASSO model with Optimal Class Weights
2  def build_op_en(X_train, y_train, X_test, threshold=0.5, best_params=None):
3
4      model = ElasticNet(random_state=rs, max_iter = 1000)
5
6      # If best parameters are provided
7      if best_params:
8          model = make_pipeline(RobustScaler(),
9                                ElasticNet(alpha = best_params_en['alpha'].iloc[0],
10                                           l1_ratio = best_params_en['l1_ratio'].iloc[0],
11                                           fit_intercept = best_params_en['fit_intercept'].iloc[0],
12                                           normalize = best_params_en['normalize'].iloc[0],
13                                           precompute = best_params_en['precompute'].iloc[0],
14                                           copy_X = best_params_en['copy_X'].iloc[0],
15                                           tol = best_params_en['tol'].iloc[0],
16                                           warm_start = best_params_en['warm_start'].iloc[0],
17                                           positive = best_params_en['positive'].iloc[0],
18                                           selection = best_params_en['selection'].iloc[0]
19                                           )
20                                )
21
22      model.fit(X_train, y_train)
23
24      # Get Prediction
25      pred = model.predict(X_test)
26      #pred = pred * (y_train.std() + y_train.mean())
27      test_ids = X_test.index
28      pred = pd.DataFrame({'Id': test_ids, 'Predicted Price': pred})
29
30      # Get Model Directory
31      print(colored("Model Directory: \n", 'cyan', attrs=['bold']))
32      + colored(dir(model), 'magenta', attrs=['bold']))
33
34      # Get Model Features
35      feature_names = model[-1].get_feature_names_out() #https://scikit-learn.org/stable/auto_examples/inspection/plot_Li
36      X_train_preprocessed = pd.DataFrame(model[-1].transform(X_train), columns=feature_names) # X_train to compute stand
37      features = pd.DataFrame(model[-1].coef_ * X_train_preprocessed.std(axis=0), columns=["Coefficients"], index=feature_
38
39  else:
40      model.fit(X_train, y_train)
41
42      # Get Prediction
43      pred = model.predict(X_test)
44      #pred = pred * (y_train.std() + y_train.mean())
45      test_ids = X_test.index
46      pred = pd.DataFrame({'Id': test_ids, 'Predicted Price': pred})
47
48      # Get Model Directory
49      print(colored("Model Directory: \n", 'cyan', attrs=['bold']))
50      + colored(dir(model), 'magenta', attrs=['bold']))
51
52      # Get Model Features
53      feature_names = model._check_feature_names
54      X_train_preprocessed = X_train.copy() #.T # Transpose X_train to compute standard deviation of the related feature
55      # Multiply model coefficients by std of features to reduce all coefficients to same unit of measure.
56      features = pd.DataFrame(model.coef_ * X_train_preprocessed.std(axis=0), columns=["Coefficients"], index=X_train.co
57
58  features['Features'] = features.index
59  #features = features.loc[features['Coefficients'] > 0] # Uncomment if only features with Coeff > 0 are required
60  features = features.sort_values(by=['Coefficients'], ascending=False).reset_index(drop=True)
61
62  # Plot Features
63  data = features.copy()
64  data = data[data["Coefficients"] != 0]
65  data.set_index('Features', inplace=True)
66  data.plot.barh(figsize=(30,10), color='green')
67
68  title = "Feature Importance in Predicting " + target
69  plt.title(title)
70  plt.axvline(x=0, color="red")
71  plt.xlabel("Coefficient Values Corrected by Features' std. dev.")
72  plt.subplots_adjust(left=0.3)
73
74  # Get Model Results
75  res = result(model)
76
77  return model, res, features, pred

```

3a) EN Model 1 Without Optimal Hyperparameter Tuning: A simple algorithm was created without any hyperparameter tuning:

Model 5: ELASTIC-NET Without 'Optimal Hypertuned Parameters'

```

1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 model5, res5, features5, pred5 = build_op_en(X_train, y_train, X_test, best_params=None) # Call Method to 'Build Logistic Re
3
4 #Collect & Append Results
5 r5 = np.array2string(res5)####
6 r5 = pd.DataFrame([r5.split(';') for x in r5.split('\n')])
7 r5 = r5.rename(columns={0: 'SCORE'})
8 r5['MODEL'] = 'Elastic-Net'
9 r5['MODEL#'] = 'Model 5'
10 r5

```

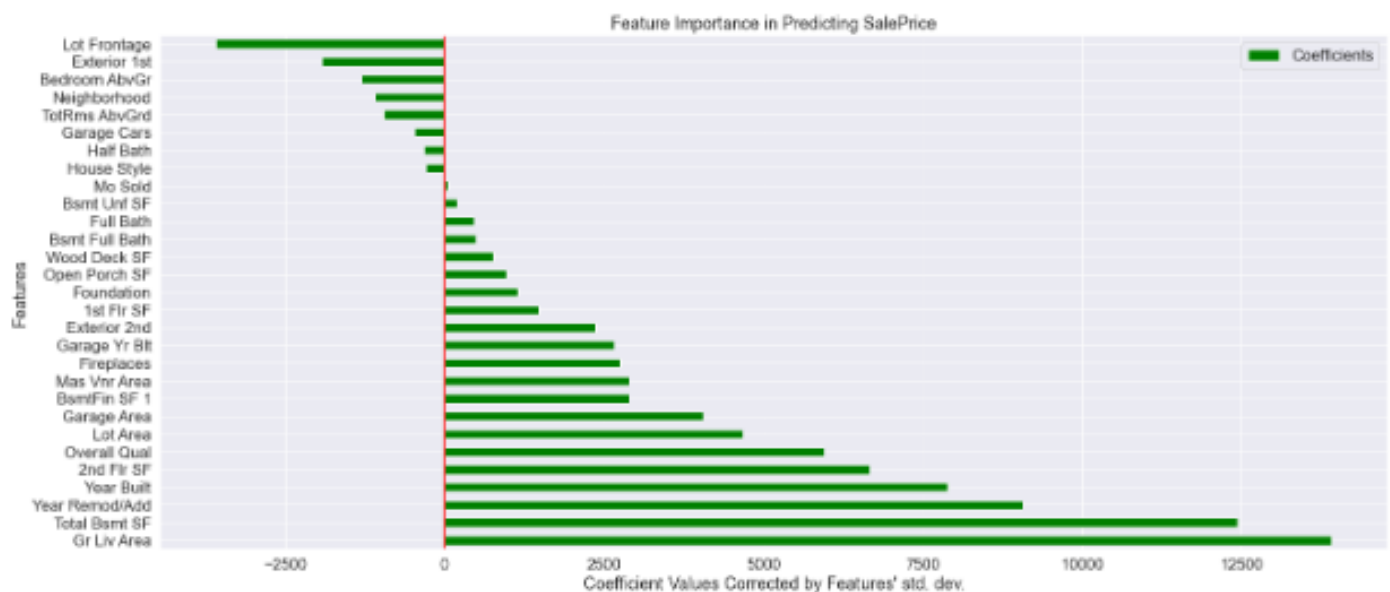
Model Directory:

```

['_abstractmethod_', '_class_', '_delattr_', '_dict_', '_dir_', '_doc_', '_eq_', '_format_', '_ge_', '_geta
ttribute_', '_getstate_', '_gt_', '_hash_', '_init_', '_init_subclass_', '_le_', '_lt_', '_module_', '_ne_',
'_new_', '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_setstate_', '_sizeof_', '_str_', '_subclasshook_',
'_weakref_', '_abc_impl', '_check_feature_names', '_check_n_features', '_decision_function', '_estimator_type', '_get_par
am_names', '_get_tags', '_more_tags', '_repr_html_', '_repr_html_inner', '_repr_mimebundle_', '_set_intercept', '_validate_dat
a', '_alpha', '_coef_', '_copy_X_', '_dual_gap_', '_feature_names_in_', '_fit', '_fit_intercept', '_get_params', '_intercept_', '_li_rati
o', '_max_iter_', '_n_features_in_', '_n_iter_', '_normalize_', '_path_', '_positive_', '_precompute_', '_predict_', '_random_state_', '_score_',
'_selection_', '_set_params_', '_sparse_coef_', '_tol_', '_warm_start_']

```

	SCORE	MODEL	MODEL#
0	-1.32680608e+09	Elastic-Net	Model 5



3b) EN Model 2 WITH Optimal Hyperparameter Tuning: A modified algorithm with auto hyper-tuned parameters was created using random search method described above:

Model 6: ELASTIC-NET With 'Optimal Hypertuned Parameters'

```

1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 model6, res6, features6, pred6 = build_op_en(X_train, y_train, X_test, best_params=optimal_params_en) # Call Method to 'Build
3
4 #Collect & Append Results
5 r6 = np.array2string(res6)####
6 r6 = pd.DataFrame([r6.split(';') for x in r6.split('\n')])
7 r6 = r6.rename(columns={0: 'SCORE'})
8 r6['MODEL'] = 'Elastic-Net Optimal'
9 r6['MODEL#'] = 'Model 6'
10 r6

```

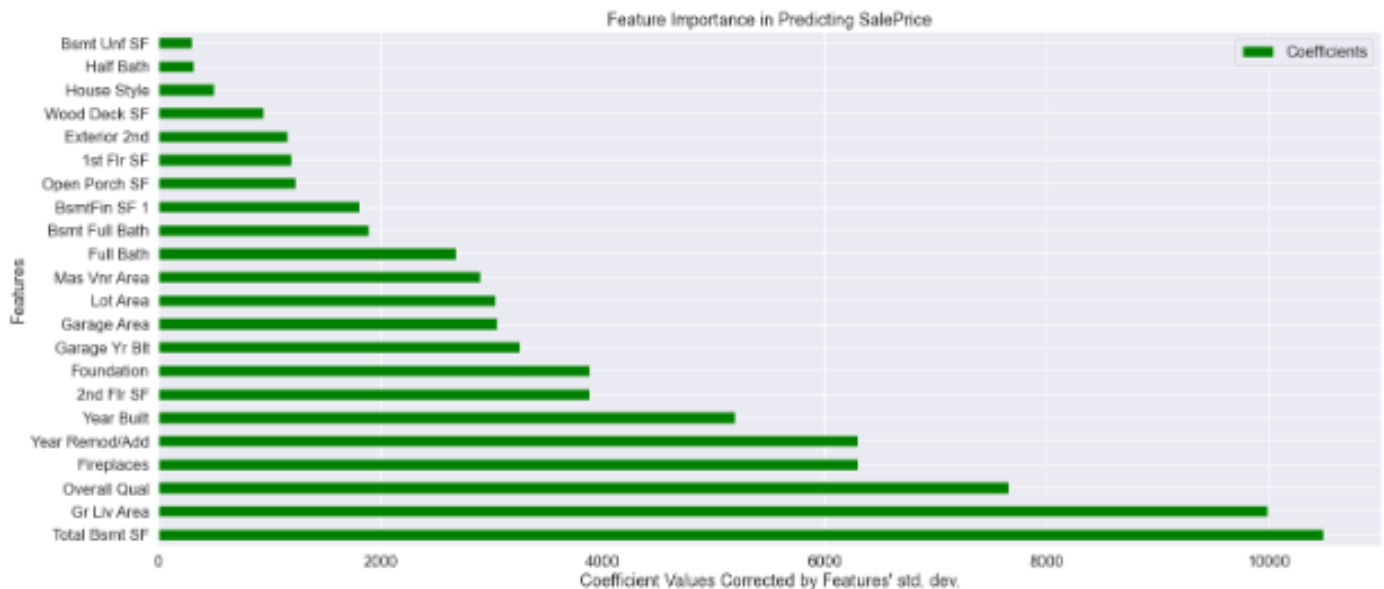
Model Directory:

```

['_abstractmethods_', '_annotations_', '_class_', '_delattr_', '_dict_', '_dir_', '_doc_', '_eq_', '_format_
_', '_ge_', '_getattr_', '_getitem_', '_getstate_', '_gt_', '_hash_', '_init_', '_init_subclass_', '_le_
_', '_len_', '_lt_', '_module_', '_ne_', '_new_', '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_setsta
te_', '_sizeof_', '_sklearn_is_fitted_', '_str_', '_subclasshook_', '_weakref_', '_abc_impl_', '_can_inverse_transfor
m_', '_can_transform_', '_check_feature_names_', '_check_fit_params_', '_check_n_features_', '_estimator_type_', '_final_estimator_',
'_fit_', '_get_param_names_', '_get_params_', '_get_tags_', '_iter_', '_log_message_', '_more_tags_', '_replace_estimator_', '_repr_htm
l_', '_repr_html_inner_', '_repr_mimebundle_', '_required_parameters_', '_set_params_', '_sk_visual_block_', '_validate_data_', '_v
alidate_names_', '_validate_steps_', '_classes_', '_decision_function_', '_feature_names_in_', '_fit_', '_fit_predict_', '_fit_transform_',
'_get_feature_names_out_', '_get_params_', '_inverse_transform_', '_memory_', '_n_features_in_', '_named_steps_', '_predict_', '_predict_log_
proba_', '_predict_proba_', '_score_', '_score_samples_', '_set_params_', '_steps_', '_transform_', '_verbose_']

```

	SCORE	MODEL	MODEL#
0	-1.32126229e+09	Elastic-Net Optimal	Model 6



4) Extreme Gradient Boosting (XGB) Regression Models

A single method was employed to run XGB regression models with and without optimal hyperparameter tuning:

Method to 'Build XGB Regression WITH & Without Optimal Hyperparameters'

```

1  # Build a LASSO model with Optimal Class Weights
2  def build_op_xgb(X_train, y_train, X_test, threshold=0.5, best_params=None):
3
4      model = GradientBoostingRegressor(random_state=rs)
5
6      # If best parameters are provided
7      if best_params:
8          model = make_pipeline(RobustScaler(),
9                                GradientBoostingRegressor(learning_rate = best_params_xgb['learning_rate'].iloc[0],
10                                                         subsample = best_params_xgb['subsample'].iloc[0],
11                                                         criterion = best_params_xgb['criterion'].iloc[0],
12                                                         max_features = best_params_xgb['max_features'].iloc[0]
13                                                         )
14                                )
15      model.fit(X_train, y_train)
16
17      # Get Prediction
18      pred = model.predict(X_test)
19      #pred = pred * (y_train.std() + y_train.mean())
20      test_ids = X_test.index
21      pred = pd.DataFrame({'Id': test_ids, 'Predicted Price': pred})
22
23      # Get Model Directory
24      print(colored("Model Directory: \n", 'cyan', attrs=['bold']))
25          + colored(dir(model), 'magenta', attrs=['bold']))
26
27      # Get Model Features
28      feature_names = model[-1].get_feature_names_out() #https://scikit-learn.org/stable/auto_examples/inspection/plot_
29      X_train_preprocessed = pd.DataFrame(model[-1].transform(X_train), columns=feature_names) # X_train to compute stand
30      features = pd.DataFrame(model[-1].feature_importances_ * X_train_preprocessed.std(axis=0), columns=["Coefficients"],
31
32  else:
33      model.fit(X_train, y_train)
34
35      # Get Prediction
36      pred = model.predict(X_test)
37      #pred = pred * (y_train.std() + y_train.mean())
38      test_ids = X_test.index
39      pred = pd.DataFrame({'Id': test_ids, 'Predicted Price': pred})
40
41      # Get Model Directory
42      print(colored("Model Directory: \n", 'cyan', attrs=['bold']))
43          + colored(dir(model), 'magenta', attrs=['bold']))
44
45      # Get Model Features
46      X_train_preprocessed = X_train.copy()
47      features = pd.DataFrame(model.feature_importances_ * X_train_preprocessed.std(axis=0), columns=["Coefficients"], in
48
49  features['Features'] = features.index
50  #features = features.loc[features['Coefficients'] > 0] # Uncomment if only features with Coeff > 0 are required
51  features = features.sort_values(by=['Coefficients'], ascending=False).reset_index(drop=True)
52
53  # Plot Features
54  data = features.copy()
55  data = data[data["Coefficients"] != 0]
56  data.set_index('Features', inplace=True)
57  data.plot.barh(figsize=(30,10), color='green')
58
59  title = "Feature Importance in Predicting " + target
60  plt.title(title)
61  plt.axvline(x=0, color="red")
62  plt.xlabel("Coefficient Values Corrected by Features' std. dev.")
63  plt.subplots_adjust(left=0.3)
64
65  # Get Model Results
66  res = result(model)
67
68  return model, res, features, pred

```


4a) XGB Regression Model 1 Without Optimal Hyperparameter Tuning: A simple algorithm was created without any hyperparameter tuning:

Model 7: XGB Without 'Optimal Hypertuned Parameters'

```

1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 model7, res7, features7, pred7 = build_op_xgb(X_train, y_train, X_test, best_params=None) # Call Method to 'Build XGB Without
3
4 #Collect & Append Results
5 r7 = np.array2string(res7)####
6 r7 = pd.DataFrame([r7.split(';') for x in r7.split('\n')])
7 r7 = r7.rename(columns={0: 'SCORE'})
8 r7['MODEL'] = 'XGB'
9 r7['MODEL#'] = 'Model 7'
10 r7

```

Model Directory:

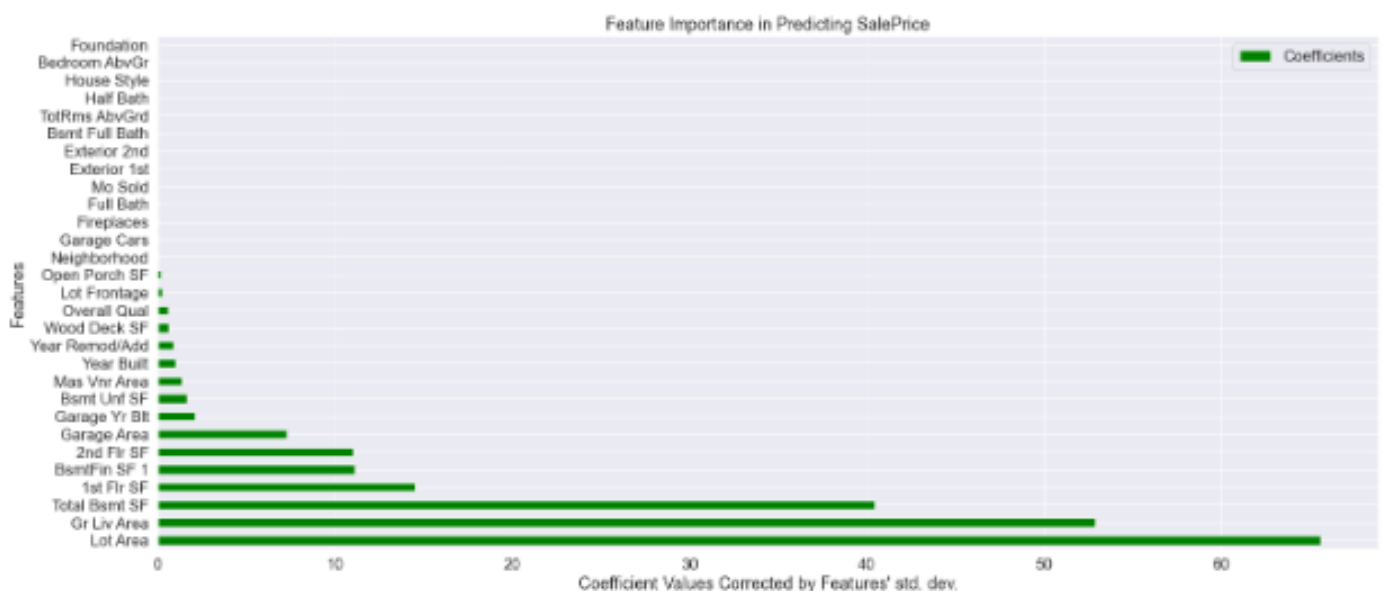
```

['_SUPPORTED_LOSS', '_abstractmethods_', '_annotations_', '_class_', '_delattr_', '_dict_', '_dir_', '_doc_', '_eq_', '_format_', '_ge_', '_getattr_', '_getitem_', '_getstate_', '_gt_', '_hash_', '_init_', '_init_subclass_', '_iter_', '_le_', '_len_', '_lt_', '_module_', '_ne_', '_new_', '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_setstate_', '_sizeof_', '_str_', '_subclasshook_', '_weakref_', '_abc_impl', '_check_feature_names', '_check_initialized', '_check_n_features', '_check_params', '_clear_state', '_compute_partial_dependence_recursion', '_estimator_type', '_fit_stage', '_fit_stages', '_get_param_names', '_get_tags', '_init_state', '_is_initialized', '_loss', '_make_estimator', '_more_tags', '_raw_predict', '_raw_predict_init', '_repr_html_', '_repr_html_inner', '_repr_mimebundle_', '_require_d_parameters', '_resize_state', '_rng', '_staged_raw_predict', '_validate_data', '_validate_estimator', '_validate_y', '_alpha', '_apply', '_ccp_alpha', '_criterion', '_estimators_', '_feature_importances_', '_feature_names_in_', '_fit', '_get_params', '_init', '_init_', '_learning_rate', '_loss_', '_loss_', '_max_depth', '_max_features', '_max_features_', '_max_leaf_nodes', '_min_impurity_decrease', '_min_samples_leaf', '_min_samples_split', '_min_weight_fraction_leaf', '_n_estimators', '_n_estimators_', '_n_features_', '_n_features_in_', '_n_iter_no_change', '_predict', '_random_state', '_score', '_set_params', '_staged_predict', '_subsample', '_tol', '_train_score_', '_validation_fraction', '_verbose', '_warm_start']

```

SCORE MODEL MODEL#

SCORE	MODEL	MODEL#
-9.9099029e+08	XGB	Model 7



4b) XGB Regression Model 2 WITH Optimal Hyperparameter Tuning: A modified algorithm with auto hyper-tuned parameters was created using random search method described above:

Model 8: XGB With 'Optimal Hypertuned Parameters'

```

1 #Running the Model...Call Method to 'Build LR Without Adjusted Class Weights'
2 model8, res8, features8, pred8 = build_op_xgb(X_train, y_train, X_test, best_params=optimal_params_xgb) # Call Method to 'Bu
3
4 #Collect & Append Results
5 r8 = np.array2string(res8)####
6 r8 = pd.DataFrame([r8.split(';') for x in r8.split('\n')])
7 r8 = r8.rename(columns={0: 'SCORE'})
8 r8['MODEL'] = 'XGB Optimal'
9 r8['MODEL#'] = 'Model 8'
10 r8

```

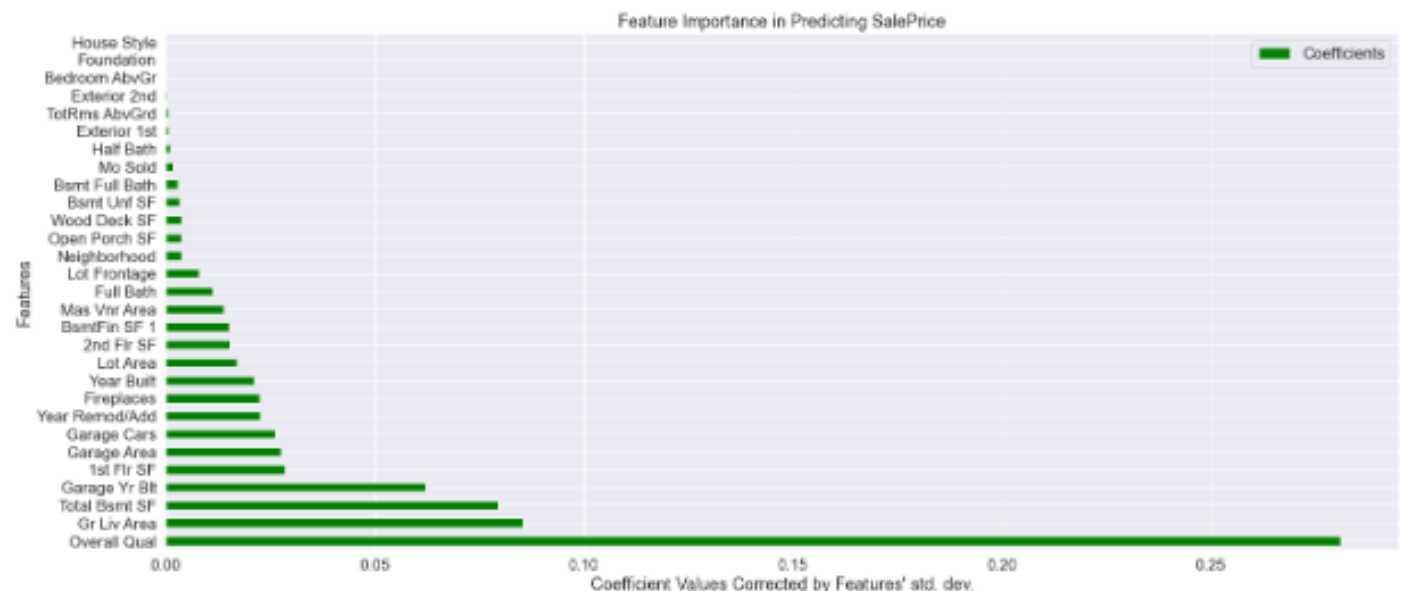
Model Directory:

```

['_abstractmethods_', '_annotations_', '_class_', '_delattr_', '_dict_', '_dir_', '_doc_', '_eq_', '_format_',
'_ge_', '_getattr_', '_getitem_', '_getstate_', '_gt_', '_hash_', '_init_', '_init_subclass_', '_le_',
'_len_', '_lt_', '_module_', '_ne_', '_new_', '_reduce_', '_reduce_ex_', '_repr_', '_setattr_', '_setsta
te_', '_sizeof_', '_sklearn_is_fitted_', '_str_', '_subclasshook_', '_weakref_', '_abc_impl_', '_can_inverse_transfor
m_', '_can_transform_', '_check_feature_names_', '_check_fit_params_', '_check_n_features_', '_estimator_type_', '_final_estimator_',
'_fit_', '_get_param_names_', '_get_params_', '_get_tags_', '_iter_', '_log_message_', '_more_tags_', '_replace_estimator_', '_repr_htm
l_', '_repr_html_inner_', '_repr_mimebundle_', '_required_parameters_', '_set_params_', '_sk_visual_block_', '_validate_data_', '_v
alidate_names_', '_validate_steps_', 'classes_', 'decision_function_', 'feature_names_in_', 'fit_', 'fit_predict_', 'fit_transform_',
'get_feature_names_out_', 'get_params_', 'inverse_transform_', 'memory_', 'n_features_in_', 'named_steps_', 'predict_', 'predict_log_
proba_', 'predict_proba_', 'score_', 'score_samples_', 'set_params_', 'steps_', 'transform_', 'verbose_']

```

	SCORE	MODEL	MODEL#
0	-9.92038691e+08	XGB Optimal	Model 8



6) Key Findings to the Main Objectives of Analysis

6a) Result Summary

Result Summary is given below:

	SCORE	MODEL	MODEL#
0	-1.3033701e+09	Lasso Optimal	Model 4
0	-1.30393778e+09	Regression Optimal	Model 2
0	-1.30784598e+09	Regression	Model 1
0	-1.30784598e+09	Lasso	Model 3
0	-1.32126229e+09	Elastic-Net Optimal	Model 6
0	-1.32680508e+09	Elastic-Net	Model 5
0	-9.9099029e+08	XGB	Model 7
0	-9.92038691e+08	XGB Optimal	Model 8

6b) Recommended Model and Justification

Recommended Model and Justification:

Model Scoring has been carried out using mean squared error (MSE).

MSE shows closeness of a regression line to a set of points and helps in finding average of a set of errors.

It takes the distances (or errors) from the points to the regression line and squares them to remove any negative signs.

Since it lends more weight to larger differences, hence, the lower the MSE, the better the forecast.

In other words, the smaller the MSE, the closer the model is to the line of best fit.

A Score of Zero would mean the model is perfect.

Therefore, when scoring a regression model with MSE, a minimal score would imply a better prediction

Hence, we will select the model with Minimal Score.

In this case, the Model Lasso Optimal is yielding Minimal Score of -1.3033701e+09.

Hence, we will select this model for Ames, Iowa House Price Prediction.

6c) Summarizing Model Drivers

Main Drivers behind top performing Lasso Optimal model are as follows:

- Uses random search method to find optimal parameters to achieve best hyper-tuning
- Runs lasso regression model with above drivers to get top contributory features

6d) Enlisting Top Contributory Factors

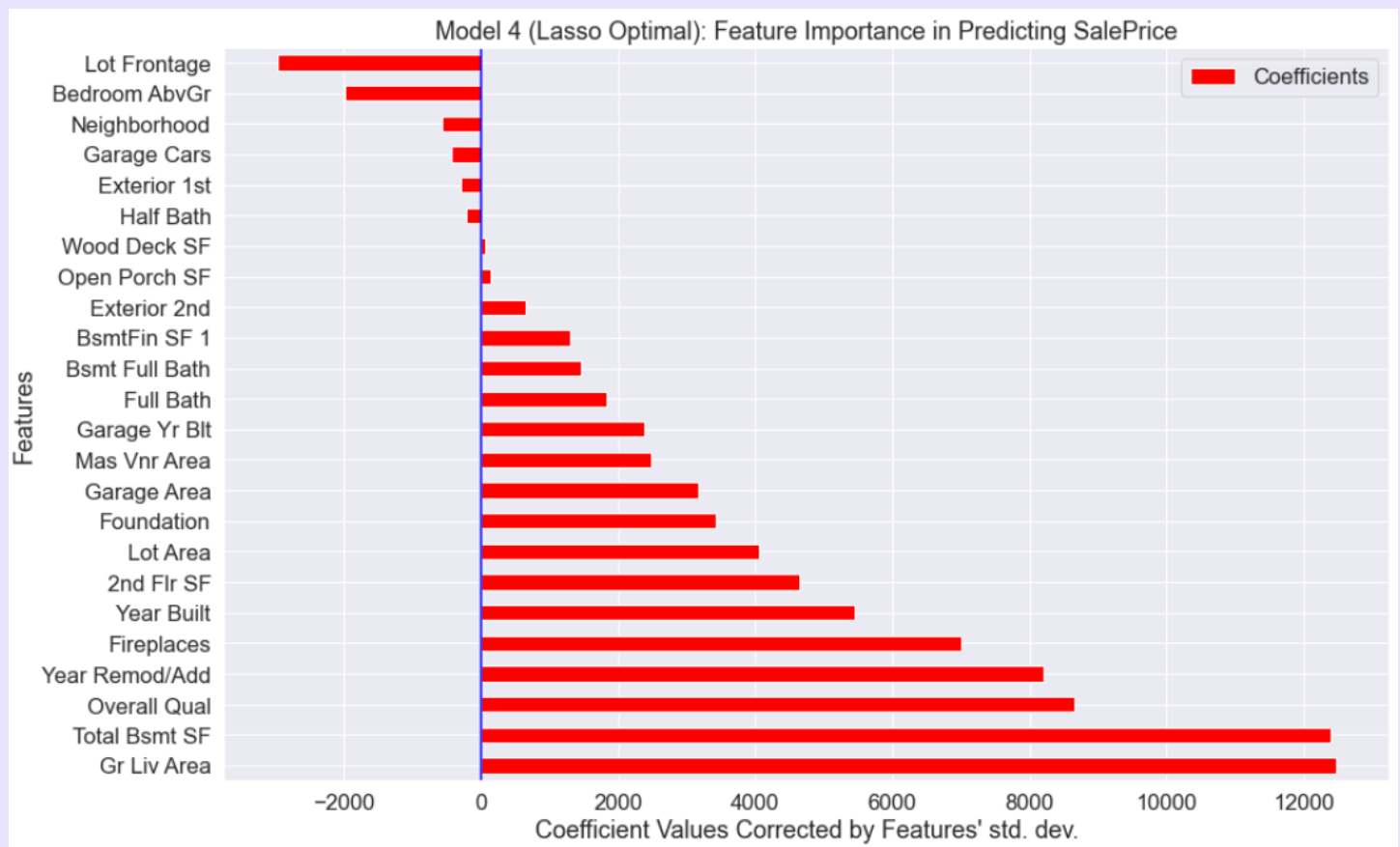
Top Factors Contributing to features driving sale value are cited under in ascending order of importance:

Contributory Features driving Sales Price:

Features Summary with Coefficient Values is given below:

	Features	Coefficients
0	Gr Liv Area	12471
1	Total Bsmt SF	12391
2	Overall Qual	8656
3	Year Remod/Add	8204
4	Fireplaces	7011
5	Year Built	5441
6	2nd Flr SF	4637
7	Lot Area	4052
8	Foundation	3413
9	Garage Area	3170
10	Mas Vnr Area	2474
11	Garage Yr Blt	2374
12	Full Bath	1830
13	Bsmt Full Bath	1448
14	BsmtFin SF 1	1302
15	Exterior 2nd	656
16	Open Porch SF	132
17	Wood Deck SF	67
39	Half Bath	-207
40	Exterior 1st	-301
41	Garage Cars	-435
42	Neighborhood	-574
43	Bedroom AbvGr	-1988
44	Lot Frontage	-2978

6e) Visualizing Top Contributory Factors Driving Sale Value



6f) Sale Price Prediction on Test/New Data

Predicted sale price for test data were extracted from top performing model as follows:

Predict Sale Price on Ames Homes' Original Test Data

```

1 # Get Target Variable Predictions
2 if model == 'Model 1':#
3     prediction = pred1.copy()
4 elif model == 'Model 2':
5     prediction = pred2.copy()
6 elif model == 'Model 3':
7     prediction = pred3.copy()
8 elif model == 'Model 4':
9     prediction = pred4.copy()
10 elif model == 'Model 5':#
11     prediction = pred5.copy()
12 elif model == 'Model 6':
13     prediction = pred6.copy()
14 elif model == 'Model 7':
15     prediction = pred7.copy()
16 elif model == 'Model 8':
17     prediction = pred8.copy()
18 else:
19     print("")
20
21 # Attach Predictions to Test Data
22 test_data = raw_data.copy()
23 #cols = ['Id']+test_data.columns.tolist()+['Predicted Price'] # Uncomment if all original columns are desired
24 cols = ['Id'] + final_list + [target] + [predicted_var] # Comment if previous line has been Uncommented
25 test_data['Id'] = test_data.index
26 test_data = test_data.merge(prediction, on='Id', how='left').sort_values(by=['Id'], ascending=True).fillna(0)
27 test_data = test_data[test_data[predicted_var] != 0]
28 test_data = test_data[cols]
29 test_data[predicted_var] = test_data[predicted_var].astype(int)
30 test_data['Price Difference'] = test_data[predicted_var] - test_data[target]
31 test_data

```

	Id	Gr Liv Area	Total Bsmt SF	Overall Qual	Remod/Add	Year	Fireplaces	Year Built	2nd Flr SF	Lot Area	Foundation	...	Wood Deck SF	Half Bath	Exterior 1st	Garage Cars	Neighborhood	Bedroom AbvGr	Frontal
0	0	1656	1080.0	6	1960		2	1960	0	31770	CBlock	...	210	0	BrkFace	2.0	NAmes	3	14
4	4	1629	928.0	5	1998		1	1997	701	13830	PConc	...	212	1	VinylSd	2.0	Gilbert	3	7
6	6	1338	1338.0	8	2001		0	2001	0	4920	PConc	...	0	0	CemntBd	2.0	StoneBr	2	4
9	9	1804	994.0	7	1999		1	1999	776	7500	PConc	...	140	1	VinylSd	2.0	Gilbert	3	6
15	15	3279	1650.0	8	2003		1	2003	1589	53504	PConc	...	503	1	CemntBd	3.0	StoneBr	4	4

These findings were compiled and subsequently exported in the form of a well formatted coloured excel report.

7) Future Directions

7a) Possible Flaws in Chosen Model

- The model uses Mean Squared Error (MSE) as scoring method which may be highly biased for higher values.
- Just like “delete_features” list, the model needs to incorporate user input list to keep features deemed necessary so that they are not automatically dropped

Outline Editable Variables

```
In [1]: 1 output_file_path = "C:/Users/fatima.s/Downloads/Ames Housing Price Prediction.xlsx"
        2 #save_plot = "C:/Users/fatima.s/Downloads/Features.png"
        3 save_plot = "Features.png"
        4 delete_features = ['Roof Matl', 'Alley'] # List Specific Features you want to retain for automated drops
```

7b) Recommendations

Following suggestions are likely to improve the model even further:

- Incorporate “keep_features” list and use it to eliminate these from final “drop” list in row 16 of project notebook

```
drop = drop1 + drop2 + drop3 + delete_features # Final List of columns to drop
```

- Introduce more models like Decision Trees and Random Forest
- Implement a method to combine best performing models to ensure enhanced performance and more effective generalization, (see, for example, 8) Useful Links, 8a-a)
- Apply other deep learning models like TensorFlow
- Because of biasness of MSE towards higher values, scoring may be substituted by Root Mean Squared Error (RMSE) which may reflect model performance whilst dealing with increased error values.
- Apply plot to depict both MSE and RMSE (see, Machines, 2022)

8) Useful Links

8a) Link to Other Useful Models

- a) <https://www.kaggle.com/code/mgmarques/houses-prices-complete-solution>
- b) <https://www.kaggle.com/code/marto24/beginners-prediction-top3>
- c) <https://www.kaggle.com/code/mchatham/ames-housing-regression>
- d) <https://www.kaggle.com/code/mkariithi/real-estate-sales-price-prediction/notebook>
- e) <https://www.kaggle.com/code/bashkeel/eda-to-ensemble-model-lasso-ridge-xgboost>
- f) <https://www.kaggle.com/code/gerlandore/advanced-house-regression-eda-model-comparison>
- g) <https://www.kaggle.com/code/prasadperera/the-boston-housing-dataset/notebook>
- h) <https://www.kaggle.com/search?q=ADVANCED+LINEAR+REGRESSION+BOSTON+HOUSE+PREDICTION>
- i) <https://www.kaggle.com/code/koki25ando/nba-salary-prediction-using-multiple-regression>

8b) Github Link to Assignment Notebook and Other Files

<https://github.com/FATIMASP/IBM-MACHINE-LEARNING-CERTIFICATION/tree/main/Supervised%20Machine%20Learning:%20Regression>

References

Herath, S. & Maier, G., 2010. *The hedonic price method in real estate and housing market research: a review of the literature, (pp. 1-21)*, Vienna, Austria: University of Economics and Business: Institute for Regional Development and Environment,

Machines, I. L., 2022. *Mean Squared Error*. [Online]

Available at: https://insidelearningmachines.com/mean_squared_error/

[Accessed 24 09 2022].