

Exploring the Paradigm Shift: Harnessing Data Analytics for Real - World Applications

D. Diana Julie M. E¹, Abinesh Mathivanan²

^{1, 2}KIT – Kalaingar Karunanidhi Institute of Technology, Department of Artificial Intelligence and Data Science, Coimbatore, India

Email: dianajuliekit[at]gmail.com

Email: abineshmathivanan31[at]gmail.com

Abstract: *Data handling and extraction are essential steps in any data science project. In Python, there are a variety of libraries that can be used to perform these tasks. This paper reviews the most popular libraries for handling and extracting datasets, including Pandas, NumPy, SciPy, Matplotlib, and scikit - learn. Pandas is a powerful library for data analysis and manipulation. It provides high - level data structures and tools for working with structured and time series data. NumPy is a library for scientific computing. It provides high - performance multidimensional arrays and a wide range of mathematical functions. SciPy is a collection of scientific computing tools that build on top of NumPy. It provides functions for linear algebra, numerical integration, and signal processing. Matplotlib is a library for creating visualizations. It provides a wide range of plotting functions for 2D and 3D data. The paper discusses the implementation of these libraries into a real world application and provides examples of how they can be used to handle and extract datasets. The paper also discusses how these libraries can be used to run simulation models of the data using ML Models. The paper concludes by discussing the future of data handling and extraction in Python, and how the development of new libraries is likely to continue to improve the tools available to data scientists.*

Keywords: Python, data handling, data extraction, Pandas, NumPy, SciPy, Matplotlib, scikit - learn, machine learning, simulation models

1. Introduction

In today's data - driven world, the ability to handle and process data effectively has become crucial for making informed decisions and driving meaningful outcomes. This paper aims to provide a comprehensive yet accessible overview of methods for handling and processing data sets using Python libraries, focusing on their practical implementation in real - life scenarios.

The presented techniques are designed to enable researchers, practitioners, and data enthusiasts to leverage the power of Python libraries in their data analysis endeavors. By exploring different methods of data handling and processing, we aim to bridge the gap between raw data and actionable insights, making the implementation of data sets in real - life applications more attainable.

Throughout the paper, we will delve into a real world application model which includes data cleaning, preprocessing, feature extraction, and transformation. We will showcase how Python libraries, such as scikit - learn, provide a wealth of tools and functions that simplify these processes. By utilizing these libraries, users can streamline the handling and processing of data sets, making them suitable for implementation in real - life scenarios.

The focus of this paper is to present the methods and techniques in a clear and accessible manner, ensuring that readers with varying levels of expertise can easily grasp and implement them. We will provide practical examples and step - by - step guidelines to facilitate understanding and promote practical application.

Overall, by providing an accessible overview of data handling and processing methods, we aim to empower

readers to unlock the full potential of their data and drive meaningful outcomes in their respective fields.

2. Libraries for Data Handling and Manipulation

Python has several widely used libraries for data handling and manipulation. Among the most popular ones are:

1) Numpy: Numpy is a Python library commonly used for numerical and scientific computing. It provides a powerful and efficient way to manipulate large multidimensional arrays and matrices. A key feature of Numpy is the N - d array (N - dimensional array) object that allows efficient storage and manipulation of homogeneous data. Numpy's array manipulation routines, such as reshaping, slicing, and indexing, enable efficient data manipulation and extraction. These operations provide a convenient way to extract subsets of data, reorder elements, and perform transformations on arrays. Numpy arrays can also be easily combined, split, and stacked, allowing for flexible data manipulation and aggregation.

Moreover, Numpy integrates seamlessly with other Python libraries and tools, such as Matplotlib for data visualization and Pandas for data analysis. Numpy arrays can be directly used as input for plotting functions and can be easily converted to Pandas Data Frames for further data exploration and analysis.

Numpy's advanced features, such as linear algebra operations, Fourier transforms, and random number generation, make it a versatile tool for scientific computing, data analysis, and numerical simulations. These features enable users to perform complex computations and analyses with ease.

Volume 12 Issue 6, June 2023

www.ijsr.net

Licensed Under Creative Commons Attribution CC BY

2) Scipy: Scipy offers a range of numerical algorithms for solving mathematical problems, including integration techniques and optimization routines. Its integration modules provide efficient implementations of numerical integration and quadrature methods, enabling accurate computation of definite and indefinite integrals. Scipy's optimization routines are useful in various applications, such as parameter estimation and model fitting.

It also includes modules for interpolation, enabling the generation of smooth curves and surfaces from sparse or irregularly spaced data points. Its interpolation methods include linear, polynomial, spline, and radial basis function interpolation, providing a variety of options for one-dimensional and multi-dimensional data.

Linear algebra operations are another strength of Scipy, with functions for solving systems of linear equations, computing matrix decompositions, and performing matrix operations like inversion, determinant computation, and eigenvalue calculations. These capabilities make Scipy well-suited for tasks involving linear algebra, such as regression analysis and principal component analysis.

Furthermore, Scipy includes a comprehensive statistics module that provides numerous statistical functions and probability distributions. It allows for statistical analysis, hypothesis testing, descriptive statistics, and other statistical computations. The integration with Numpy arrays and other Scipy modules facilitates seamless data handling and analysis within a consistent scientific computing environment.

3) Matplotlib: Matplotlib is a popular Python library often used for visualization and graph creation. It provides a flexible and comprehensive set of functions for generating a variety of static, animated and interactive plots, diagrams and figures.

It allows you to create line charts, scatter plots, bar charts, histograms, pie charts, heat maps and many other types of visualizations. This library provides fine-grained control over many aspects of plotting, including axes, labels, titles, colors, line styles, markers, annotations, and more. This allows for precise tuning and design of visualizations to your specific requirements. It supports a variety of output formats, allowing users to save plots as image files (PNG, JPEG, SVG, etc.) or directly embed them in various file formats (PDF, LaTeX, etc.). This makes it ideal for producing publication-quality figures in research papers, reports, presentations, websites, etc.

In addition to its core functionality, Matplotlib seamlessly integrates with other libraries such as NumPy and Pandas to enable efficient data processing and manipulation. It can directly represent data from NumPy arrays, pandas dataframes, and other data structures, simplifying the process of visualizing data stored in these formats.

Matplotlib is well documented and offers an extensive collection of examples, tutorials and documentation. This allows users to easily get started, learn the capabilities of the library, and find solutions to specific plotting challenges.

4) Scikit - learn: Scikit - learn is a widely adopted and highly regarded Python library that provides a rich set of functionalities for machine learning tasks, data analysis, and predictive modeling. With its intuitive and consistent API, extensive algorithm collection, and efficient implementation, scikit - learn has become a go - to choice for researchers and practitioners in the field of machine learning.

One of the key strengths of scikit - learn lies in its comprehensive support for various stages of the machine learning pipeline. It offers a wide range of tools for data preprocessing, feature extraction, feature selection, and feature scaling, allowing users to effectively handle and manipulate datasets of different sizes and complexities. Additionally, scikit - learn provides robust implementations of diverse machine learning algorithms, including both supervised and unsupervised approaches, making it suitable for a broad range of applications. It seamlessly integrates with other popular Python libraries such as NumPy, SciPy, and Matplotlib, allowing for seamless data exchange, visualization, and advanced scientific computing capabilities.

The versatility of scikit - learn extends beyond traditional machine learning tasks. It also supports more advanced techniques such as model evaluation, parameter tuning, and ensemble learning, facilitating the development of robust and accurate predictive models. Additionally, scikit - learn offers utilities for working with multi-class classification, text data, and imbalanced datasets, addressing various challenges encountered in real-world scenarios.

5) Seaborn: Seaborn is a popular Python library that enhances the capabilities of data visualization, building upon the functionality provided by Matplotlib. With its intuitive and user-friendly interface, Seaborn simplifies the process of creating visually appealing and informative statistical graphics. It offers a diverse range of built-in themes, color palettes, and specialized plot types, allowing researchers and data scientists to effectively communicate complex information through visual representations.

One of the key advantages of Seaborn is its seamless integration with Pandas, a powerful data manipulation library. This integration enables users to effortlessly handle and analyze datasets, making data exploration and analysis more efficient. Moreover, Seaborn provides extensive support for multi-dimensional datasets, allowing researchers to uncover intricate relationships and patterns within their data.

The rich set of visualization options offered by Seaborn empowers researchers to generate compelling visualizations that effectively convey insights and facilitate data-driven decision-making. Whether it's exploring data distributions, examining relationships between variables, or visualizing statistical models, Seaborn provides a versatile toolkit that enhances the visual exploration and understanding of data.

3. Real world Implementation with ML and Data Analysis

The below code is focused on performing Exploratory analysis on House Prices: Advanced Regression Techniques. You could find the dataset for the Housing Prices in this link <https://www.kaggle.com/competitions/house-prices-advanced-regression-techniques/data>

We could start the analysis by importing the necessary libraries for Data Analysis and manipulation of the housing price:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats
from scipy.stats import norm, skew
from sklearn.preprocessing import StandardScaler
import warnings
warnings.filterwarnings('ignore')
```

```
sns.set_style("darkgrid")
plt.style.use("fivethirtyeight")
pd.pandas.set_option('display.max_columns', None)
%matplotlib inline
```

The above snippet sets the visual styles for plots using the "darkgrid" style for Seaborn and the "fivethirtyeight" style for Matplotlib. It also configures Pandas to display all columns in DataFrames and enables inline plotting in Jupyter Notebook.

```
train = pd.read_csv("train.csv")
test = pd.read_csv("test.csv")
train.shape
train.columns
```

The above snippet reads, "train.csv" and "test.csv", and stores the data in separate DataFrames "train" and "test". The **train.shape** line returns the dimensions of the "train" DataFrame, while **train.columns** returns the column names of the "train" DataFrame, **train.head()** displays a preview of the first few rows of the "train" DataFrame, providing a quick glimpse into the data's structure and values.

Output:

(1460, 81)

```
Index (['Id', 'MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street', 'Alley', 'LotShape', 'Land - Contour', 'Utilities', 'LotConfig', 'LandSlope', 'Neighborhood', 'Condition1', 'Condition2', 'BldgType', 'HouseStyle', 'OverallQual', 'OverallCond', 'YearBuilt', 'YearRemodAdd', 'RoofStyle', 'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'MasVnrArea', 'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinSF1', 'BsmtFinType2', 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', 'Heating', 'HeatingQC', 'CentralAir', 'Electrical', '1stFlrSF', '2ndFlrSF', 'LowQualFinSF', 'GrLivArea', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath',
```

```
'HalfBath', 'BedroomAbvGr', 'KitchenAbvGr', 'KitchenQual', 'TotRmsAbvGrd', 'Functional', 'Fireplaces', 'FireplaceQu', 'GarageType', 'GarageYrBlt', 'GarageFinish', 'GarageCars', 'GarageArea', 'GarageQual', 'GarageCond', 'PavedDrive', 'WoodDeckSF', 'OpenPorchSF', 'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'PoolArea', 'PoolQC', 'Fence', 'MiscFeature', 'MiscVal', 'MoSold', 'YrSold', 'SaleType', 'SaleCondition', 'SalePrice'], dtype='object')
```

Out of the 81 features, some may not contribute significantly to the prediction. After careful evaluation, it is observed that one additional feature **TotRmsAbvGrd**, could also play a crucial role in predicting the SalePrice. However, it is noted that the **Neighborhood** feature does not provide substantial information in this context. The key features to focus on are **OverallQual**, **YearBuilt**, **TotalBsmtSF**, **GrLivArea**, and **TotRmsAbvGrd**.

```
train["SalePrice"].describe()
```

Output:

```
count 1460.000000
mean 180921.195890
std 79442.502883
min 34900.000000
25% 129975.000000
50% 163000.000000
75% 214000.000000
max 755000.000000
Name: SalePrice, dtype: float64
```

Based on the statistical summary of the feature, we can observe a noticeable skewness without the need for visualizing the histogram or KDE plot. The discrepancy between the mean and median values (25%) indicates the presence of skewness. Additionally, while 75% of the price values fall within the range of 35000 to 215000, the maximum value significantly deviates from this range. To validate our observation, we will plot the data and confirm our conclusion,

```
= sns.distplot(train.SalePrice)
```

Output:

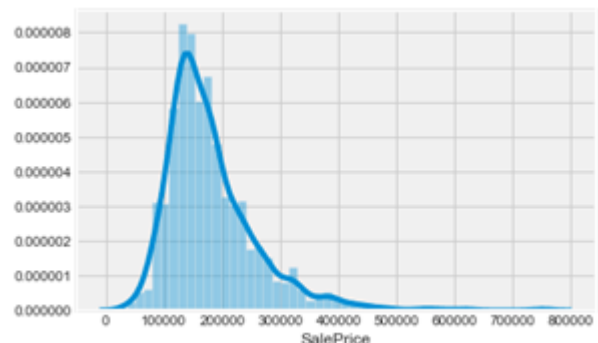


Figure 1: Simulation of Housing Sale Prices using Matplotlib

Our initial observation was accurate as the data exhibits a right-skewed distribution, with the tail extending towards the right-hand side. To quantify the extent of skewness, we

will further investigate the data,
`print(f" Skewness: {train['SalePrice'].skew() }")`
`print(f" Kurtosis: {train['SalePrice'].kurt() }")`

Output:

Skewness: 1.8828757597682129

Kurtosis: 6.536281860064529

Relationship between Selected Features

Relation with Numerical Variables

The below code snippet shows the scatter plot on the linear relationship between **SalePrice** and **GrLivArea** with a slight exponential relationship between **TotalBsmtSF**

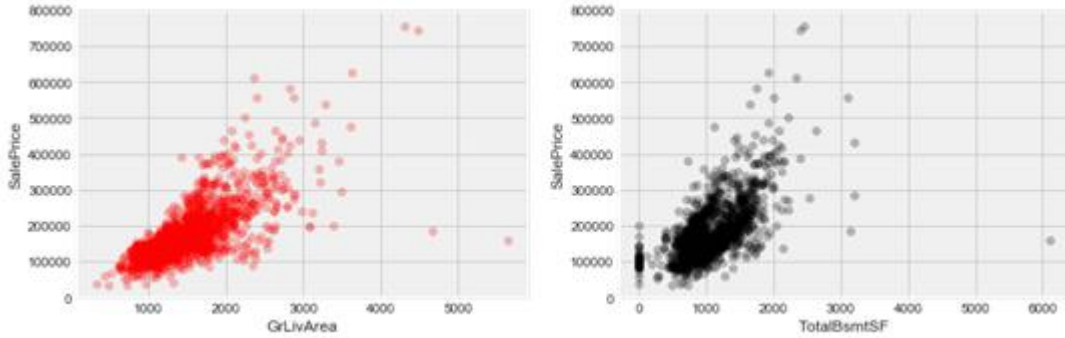


Figure 2: Scatter plot comparison between **SalePrice** and **GrLivArea**

Relation with Categorical Variables

The below snippet generates a boxplot relationship between the 'OverallQual' feature and the corresponding 'SalePrice' in the 'train' dataset. The boxplot is created by specifying

'OverallQual' as the x - axis variable and 'SalePrice' as the y - axis variable. The plot is then displayed with a restricted y - axis range of 0 to 800, 000 and it confirms that higher the quality, higher the price.

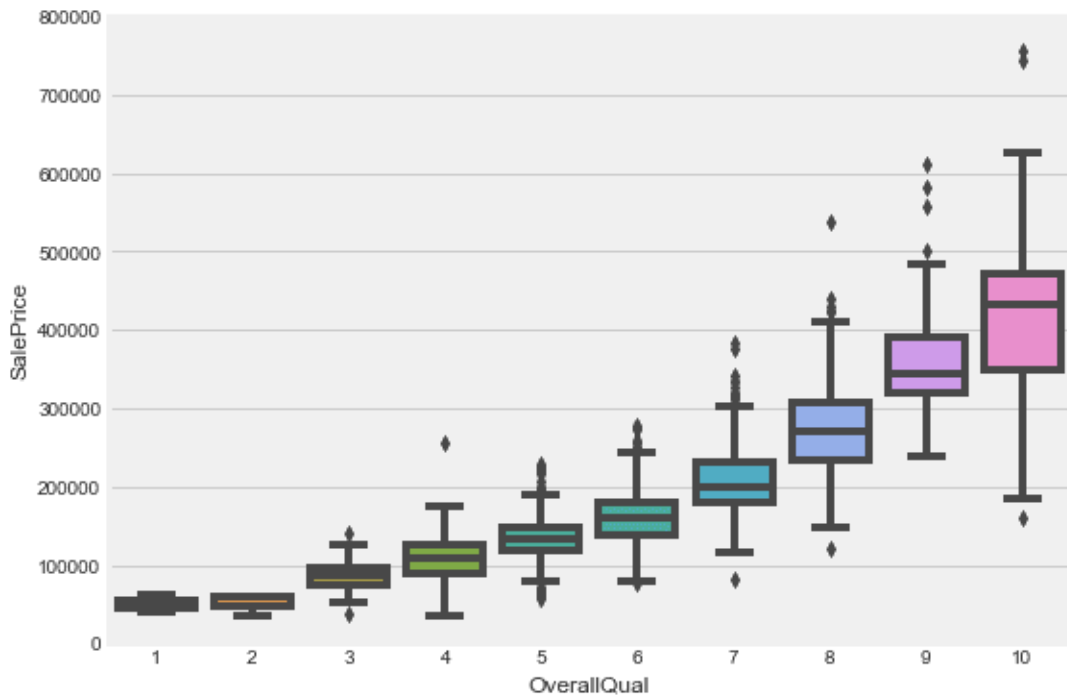


Figure 3: Boxplot comparison between **OverallQual** and **SalePrice**.

```
plt.figure(figsize=(20, 8))
var='YearBuilt'
data=pd.concat([train['SalePrice'],train[var]],axis
=1)fig=sns.boxplot(x=var,y='SalePrice',data=data)fig.axis(ymin=0,ymax=800000)
=plt.xticks(rotation=90)
```

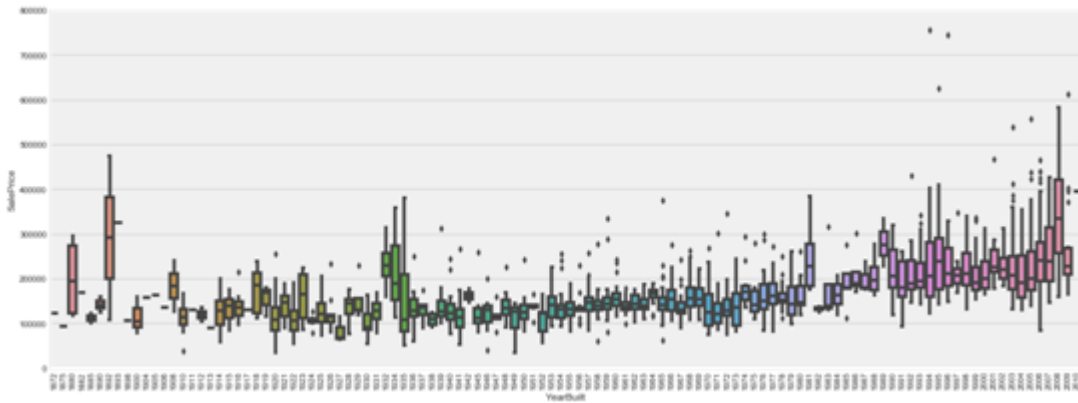


Figure 4: Boxplot comparison between **YearBuilt** and **SalePrice**. `plt.figure(figsize=(20, 6))`

```
var = 'TotRmsAbvGrd' plt.subplot(1, 2, 1)
data = pd.concat([train['SalePrice'], train[
    var]], axis=1)
= sns.boxplot(x= var, y='SalePrice', data=
    data) plt.axis(ymin=0, ymax=800000)
plt.show()
```

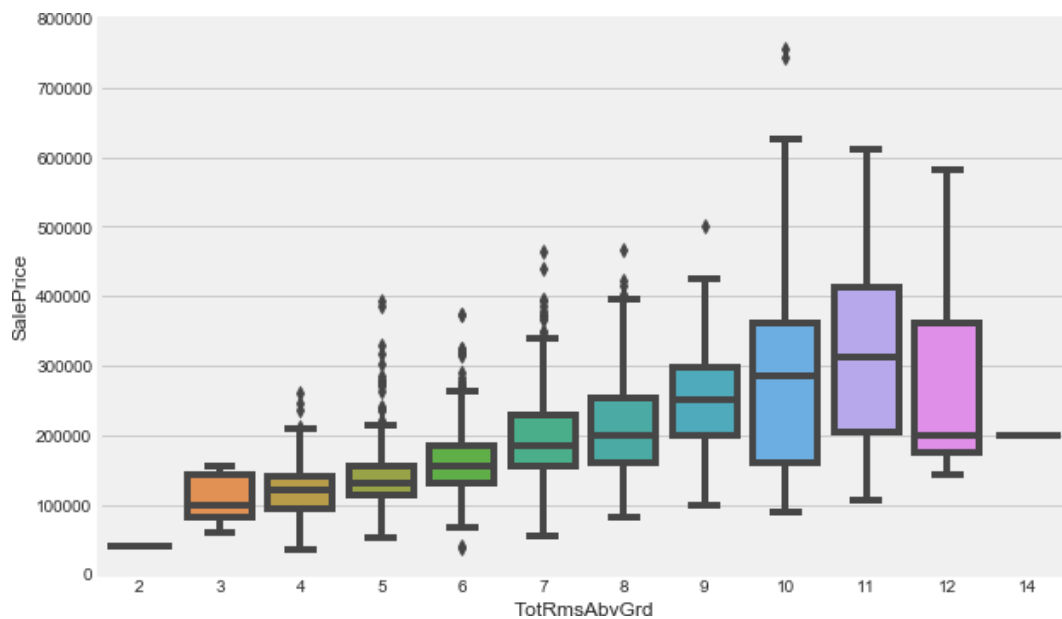


Figure 5: Boxplot comparison between **TotRmsAbvGrd** and **SalePrice**.

The variables **'GrLivArea'** and **'TotalBsmtSF'** show a positive linear relationship with **'SalePrice'**. Similarly, **'OverallQual'**, **'TotRmsAbvGrd'**, and **'YearBuilt'** also exhibit some correlation with **'SalePrice'**. Among these, the relationship is particularly strong for **'OverallQual'** as indicated by the box plot, showing an increase in sales prices with higher overall quality. The relationship is less pronounced for **'TotRmsAbvGrd'** compared to **'OverallQual'**.

Analysis on Other Variables

```
Let's have a look at all the available data subjectively,
plt.figure(figsize=(18, 10))
corrmat=train.drop('Id', 1).corr()
= sns.heatmap(corrmat, vmax=1.0, square=True, fmt
    ='.2f',
    cmap='coolwarm', annot_kws={'size':8})
plt.show()
```

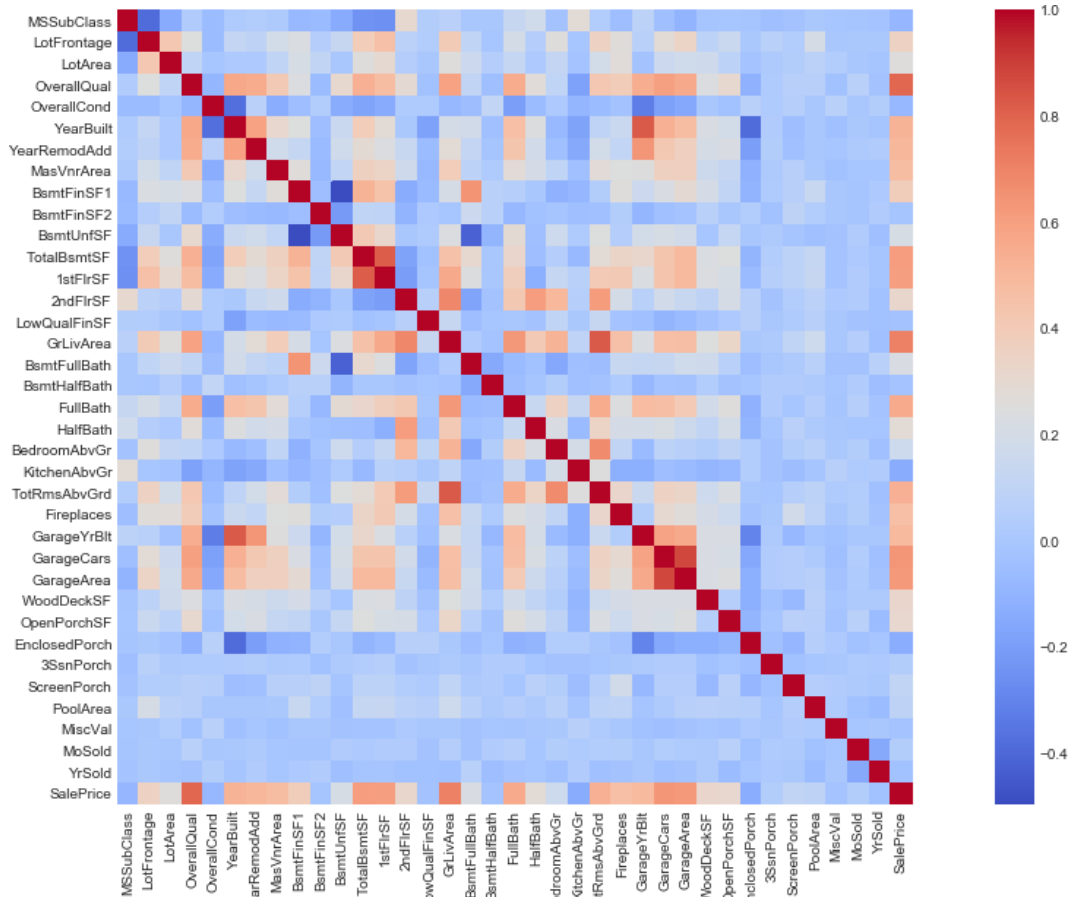


Figure 6: Heatmap visualization of the correlation matrix

As we can see, the above snippet generates a heatmap visualization of the correlation matrix of the dataset. The correlation matrix is calculated by excluding the 'Id' column from the 'train' dataset using the drop () function.

The heatmap reveals two notable correlations. The first is between 'TotalBsmtSF' and '1stFlrSF', indicating a strong relationship where both variables provide similar information, potentially leading to multicollinearity. The second correlation is observed among the 'GarageX' variables. Additionally, we can observe significant correlations with 'SalePrice' for 'GrLivArea', 'TotalBsmtSF', and 'OverallQual',

```
plt.figure(figsize=(6, 6))
```

```
k = 10 # number of variables for heatmap
cols = corrmat.nlargest(k, 'SalePrice')['SalePrice'].index
cm = np.corrcoef(train[cols].values.T)
sns.set(fontscale=1.25)
hm = sns.heatmap(cm, cbar=True, annot=True, square=True,
                 fmt='.2f', cmap='coolwarm', annot_kws={'size': 10},
                 yticklabels=cols.values, xticklabels=cols.values)
plt.show()
```

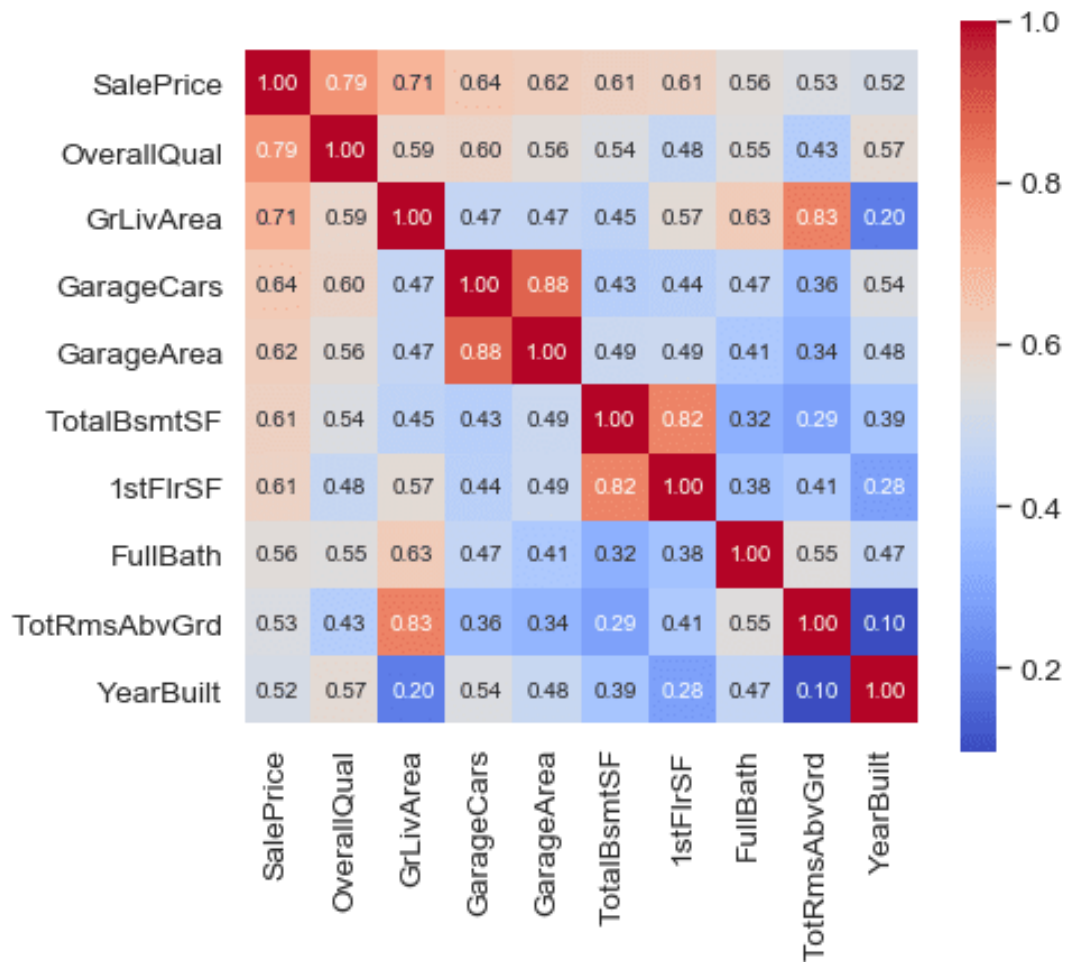


Figure 7: Heatmap of the correlation coefficients between the top 10 variables

The heatmap serves as a valuable tool for identifying such correlations and assisting in feature selection processes. It also highlights the relevance of other variables that should be considered in our analysis.

```
plt.figure(figsize=(10, 8))
sns.set()
cols = ['SalePrice', 'OverallQual', 'GrLivArea',
```

```
'GarageCars', 'TotalBsmtSF', 'FullBath', 'YearBuilt']
sns.pairplot(train[cols], size=2.5, diag_kde='kde')
plt.show()
```

The above snippet produces a pair plot that includes variables such as 'SalePrice', 'OverallQual', 'GrLivArea', 'GarageCars', 'TotalBsmtSF', 'FullBath', and 'YearBuilt' and shows the relationship between them.

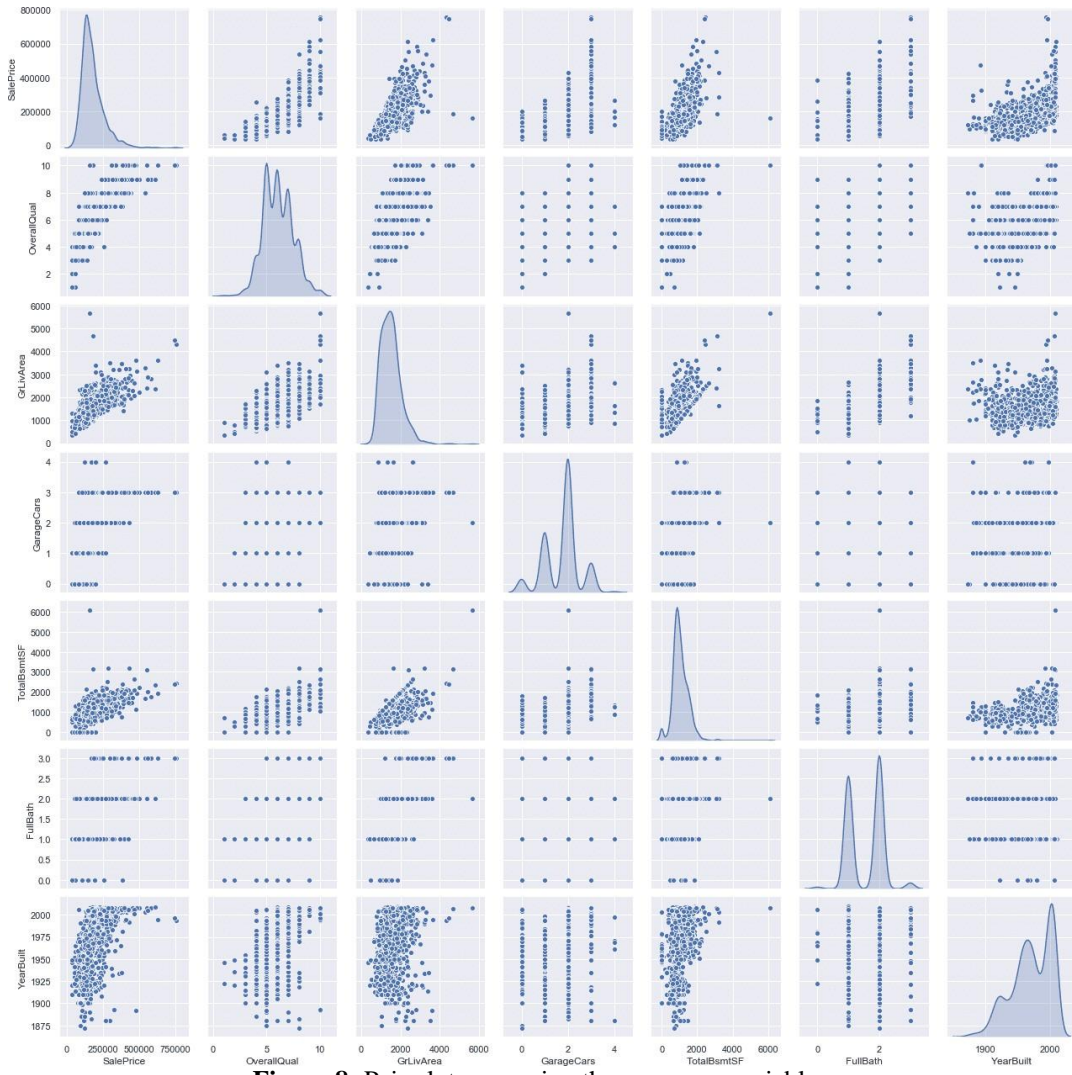


Figure 8: Pair plot comparing the necessary variables

Missing Variables

```
featureswithna = [featuresforfeaturesint
rain.columnsiftrain[features].isnull().sum()
>= 1]
```

```
a = pd.DataFrame (
'features': featureswithna,
'Total': [train[i].isnull().sum() for i in featu
reswithna], 'Missing PCT': [np.round (train [i]. is
null().sum () /
train.shape [0], 4) for i in featureswithna]
).sortvalues (by = 'Missing PCT', ascending = F
alse).resetindex (drop = True
)a.style.backgroundgradient (cmap = 'Reds')
```

We've now created a DataFrame 'a' that summarizes features in the 'train' dataset with missing values which calculates the total count and percentage of missing values for each feature and sorts them in descending order. The resulting DataFrame is styled with a red color gradient, indicating the extent of missing values in each feature.

features	Total	Missing_PCT
0 PoolQC	1453	0.995200
1 MiscFeature	1406	0.963000
2 Alley	1369	0.937700
3 Fence	1179	0.807500
4 FireplaceQu	690	0.472600
5 LotFrontage	259	0.177400
6 GarageType	81	0.055500
7 GarageYrBlt	81	0.055500
8 GarageFinish	81	0.055500
9 GarageQual	81	0.055500
10 GarageCond	81	0.055500
11 BsmtExposure	38	0.026000
12 BsmtFinType2	38	0.026000
13 BsmtFinType1	37	0.025300
14 BsmtCond	37	0.025300
15 BsmtQual	37	0.025300
16 MasVnrArea	8	0.005500
17 MasVnrType	8	0.005500
18 Electrical	1	0.000700

Figure 9: Table showing the percentage of missing values

We have several options for handling features with missing values. Features with more than 50% missing values can be dropped since they may not provide significant information. For the 'GarageX' and 'BSMTX' variables, since they are

related to garages and basements respectively, and we already have important information captured by other variables, we can consider dropping or filling them with appropriate values. Similarly, 'MasVnrArea' and 'MasVnrType' can be considered non - essential due to their correlation with 'YearBuilt' and 'OverallQual'. Lastly, for the 'Electrical' feature with just one missing value, we can either remove the corresponding row or fill it with the mode value.

In this dataset, we're replacing the missing values of the above variables with, "None", "Not available" and "No garage".

A grid of bar plots are generated where each bar plot corresponds to a feature that contains missing values in the dataset. The plots illustrate the median sale price for observations where the feature information is present versus those where it is missing. The presence or absence of a specific feature value can be visually compared to assess its

impact on the sale price. The bar plots are accompanied by informative titles and properly labeled axes, facilitating clear interpretation of the results. The resulting visualizations serve as a valuable tool for understanding the relationship between missing data in certain features and their influence on the sale price.

```
plt.figure(figsize=(20, 20))
for i, feature in enumerate(features_with_na, 1):
    plt.subplot(5, 5, i)
    data = train.copy()
    data[feature] = np.where(data[feature].isnull(), 1, 0)
    temp = data.groupby(feature)['SalePrice'].median()
    = sns.barplot(x=temp.index, y=temp.values)
    plt.title(feature)
    plt.xlabel("")
    plt.xticks([0, 1], ["Present", "Missing"])
    plt.ylabel("Sales Price", rotation=90)
    plt.tight_layout(hpad=2, wpad=2)
plt.show()
```

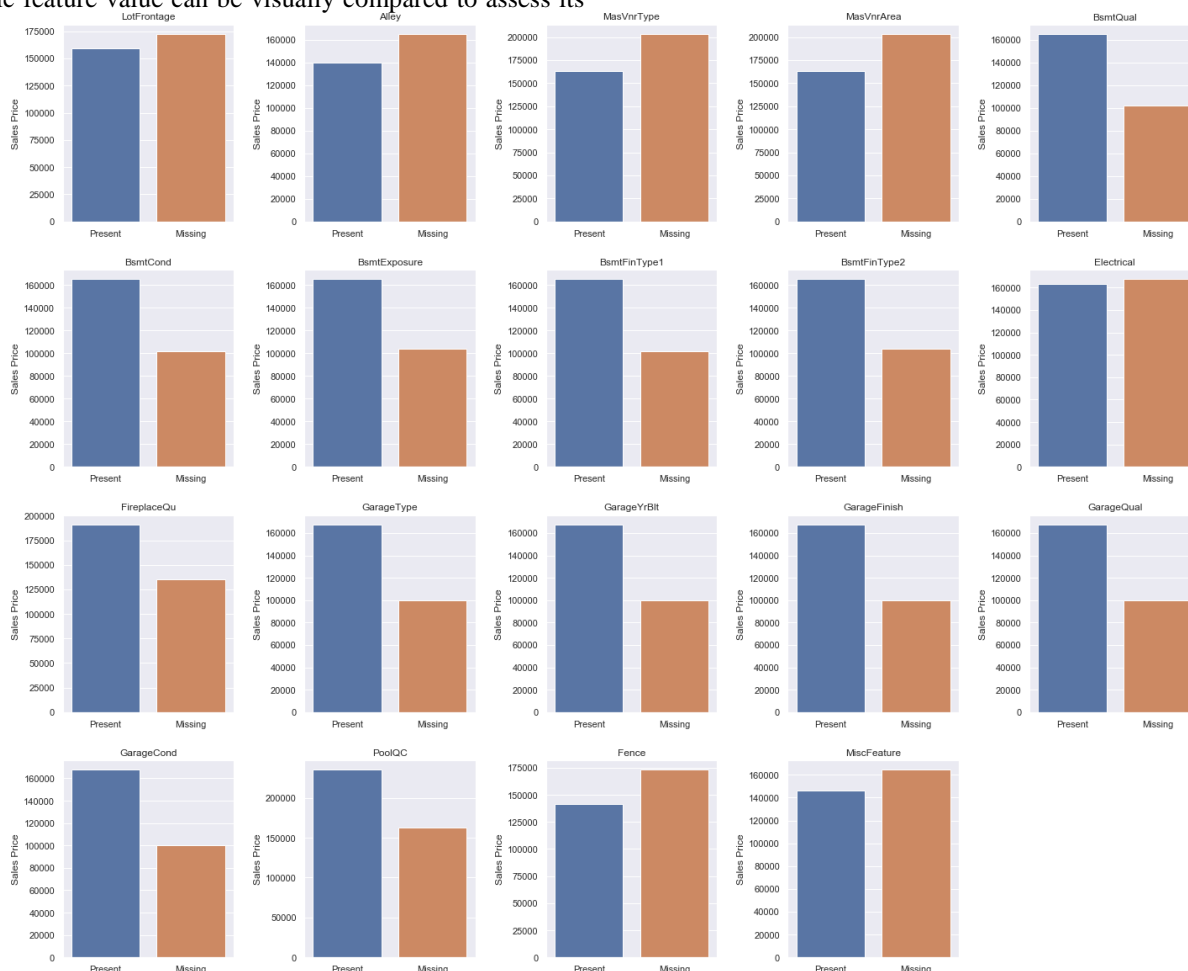


Figure 10: Grid plot simulation

Numerical Variables

Discrete Numerical Values:

Discrete numerical values are specific numeric data points that are separate and distinct, typically represented by whole numbers or integers. They are countable and finite, lacking intermediate values or decimals.

discrete features

Output:

['MSSubClass', 'OverallQual', 'OverallCond', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath', 'HalfBath', 'Bed - roomAbvGr', 'KitchenAbvGr', 'TotRmsAbvGrd', 'Fireplaces', 'GarageCars', 'PoolArea', 'MoSold']

Now, we've written a code that creates a grid of subplots to visualize the relationship between discrete features and the sale price. It iterates through each feature in the list of

discrete features and plots a bar chart showing the median sale price for each category of the feature. The data used for plotting is a copy of the training dataset. The resulting plot

allows for a comparison of how different discrete features are associated with the sale price, providing insights into their impact on the target variable.

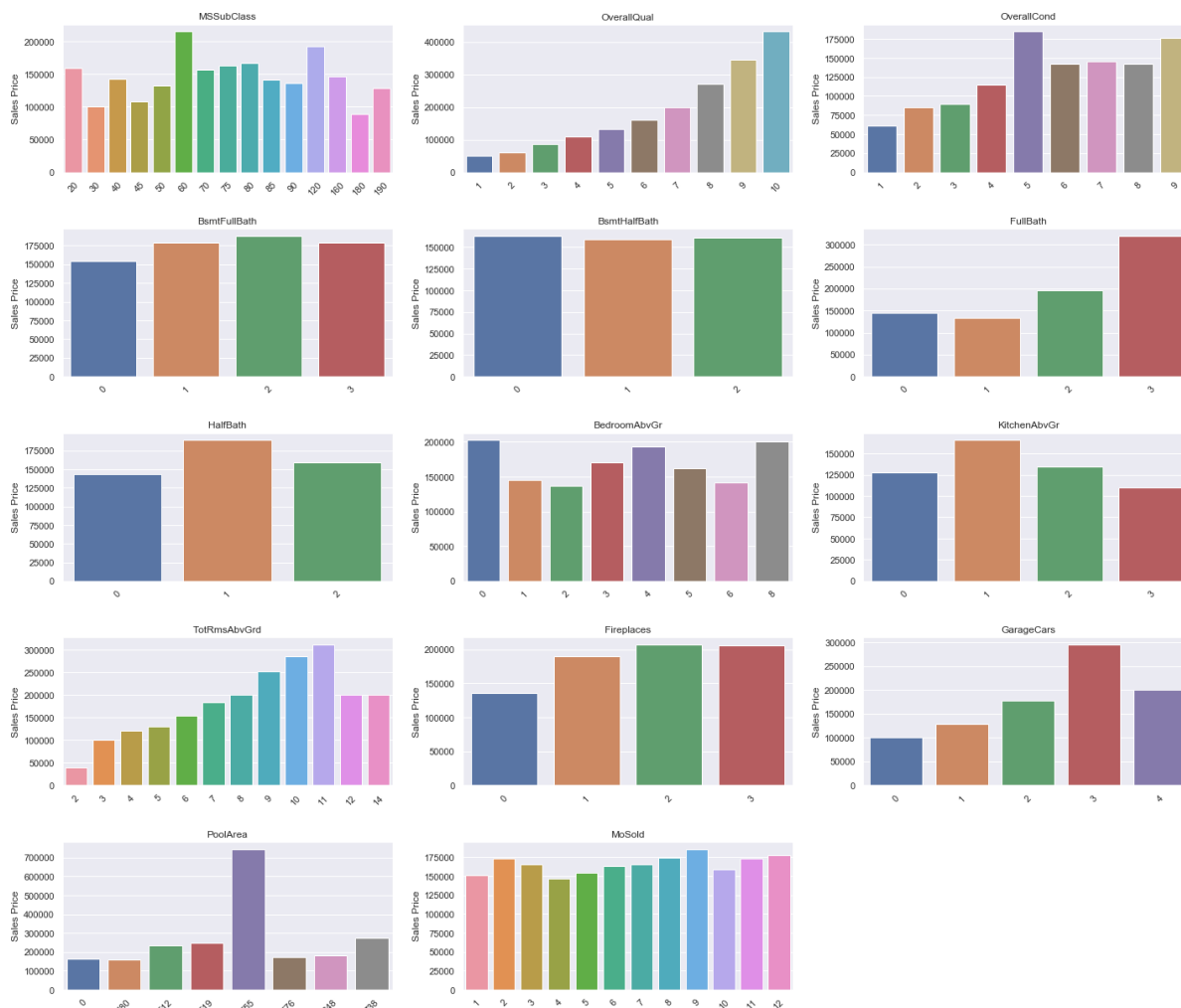


Figure 11: Relationship between Discrete Values and Sales Price The discrete numerical features exhibit varying effects on the SalePrice:

- Some features show a positive correlation, where an increase in feature value corresponds to higher SalePrice.
- Certain features have specific values that result in significantly higher or lower SalePrice.
- Several features have minimal variation in SalePrice across different feature values.
- Overall, the impact of most features on SalePrice is not immediately apparent, highlighting the need for a feature selection process to identify the most influential factors.

Continuous Numerical Values

Continuous numerical values in a dataset represent measurements or variables that can take on any value within a specific range. Unlike discrete numerical values, which have distinct categories or intervals, continuous numerical values can vary infinitely and exhibit a smooth progression.

```
print(f"Number of Continuous numerical feature: {len(continuous_num_features)}")
continuous_num_features
```

Output:

Number of Continuous numerical feature: 19

```
['LotFrontage', 'LotArea', 'MasVnrArea', 'BsmtFinSF1', 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', '1stFlrSF', '2ndFlrSF', 'LowQualFinSF', 'GrLivArea', 'GarageArea', 'WoodDeckSF', 'OpenPorchSF', 'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'MiscVal', 'SalePrice']
```

To visualize the distributions of continuous numerical features in the dataset: plt.figure(figsize=(20, 20))

```
for i, feature in enumerate(continuous_num_features, 1):
    plt.subplot(5, 4, i)
    sns.distplot(train[feature], kde=False, rug=True)
    plt.tight_layout(hpad=2, wpad=2)
plt.show()
```

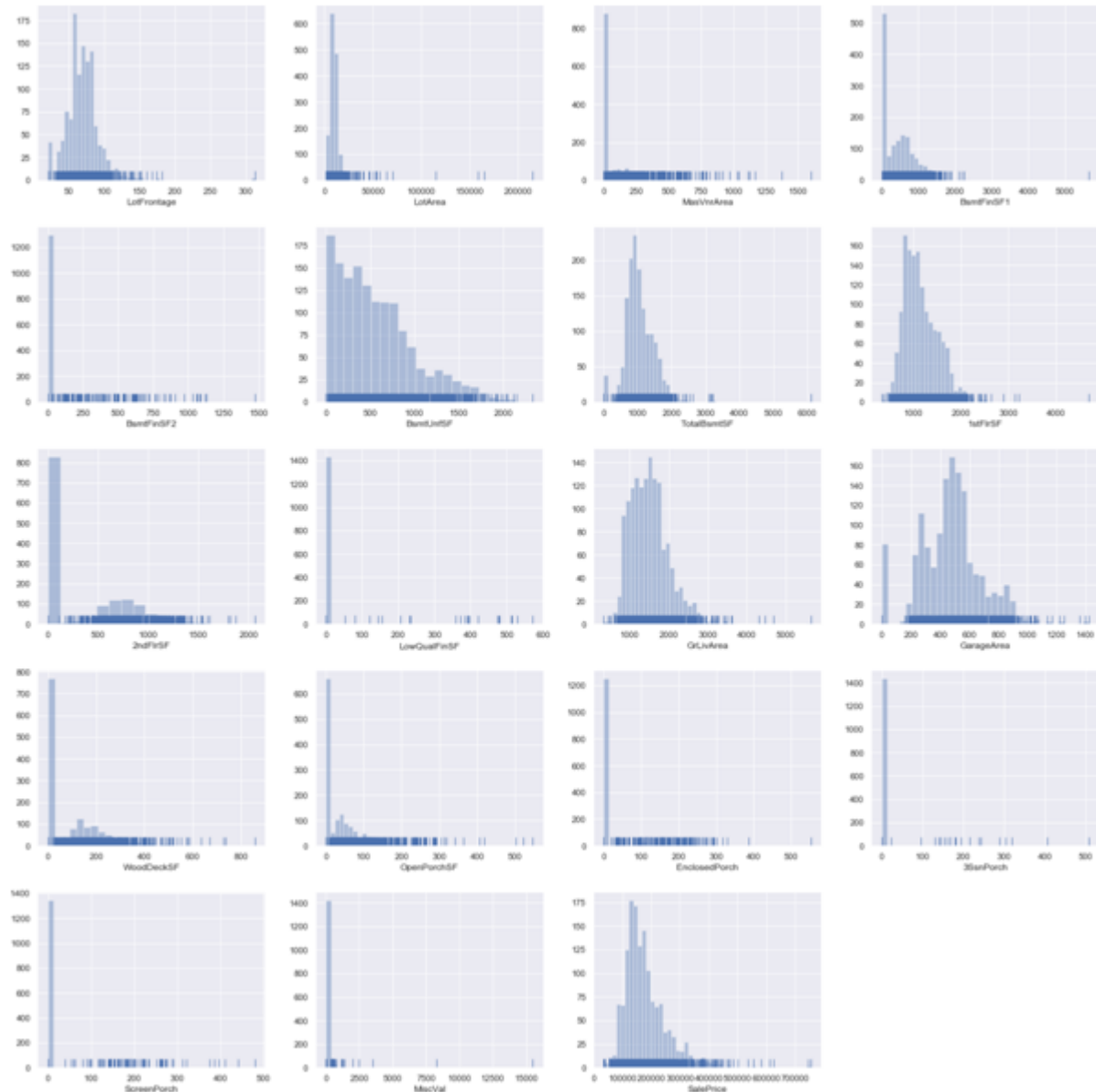


Figure 12: Simulation of Continuous Numerical Variables

The generated subplots reveal that the majority of continuous features do not exhibit a normal distribution.

Additionally, several features display significant skewness, indicating the need for data transformation to address this issue.

Logarithmic Transformation

Now we’ve to apply a log transformation to the continuous features to plot their correlation with the logarithm of **SalePrice**, and to visualize their distribution after transformation. This analysis helps to evaluate the relationship between the transformed features and the target variable, as well as the effectiveness of the log transformation in achieving a more normal distribution.

```
data = train.copy()
```

```
saleprice = np.log(train['SalePrice'])
```

```
for i, feature in enumerate(continuousnumfeatures[:-1], 1)
```

```
: data = train[feature].copy()
```

```
if 0 in data.unique(): # as log 0 is undefined
    continue
```

```
else:
```

```
data = np.log(data)
data.name = feature
= plt.figure(figsize=(15, 8))
plt.subplot(1, 2, 1)
```

```
sns.regplot(data, saleprice, fitreg=True, scatterkws="color": "black",
linekws="color": "red").set_title(f"Correlation: data.corr(saleprice)")
```

```
plt.subplot(1, 2, 2)
```

```
sns.distplot(data).set_title(f"Log transformation of: {feature}")
plt.show()
```

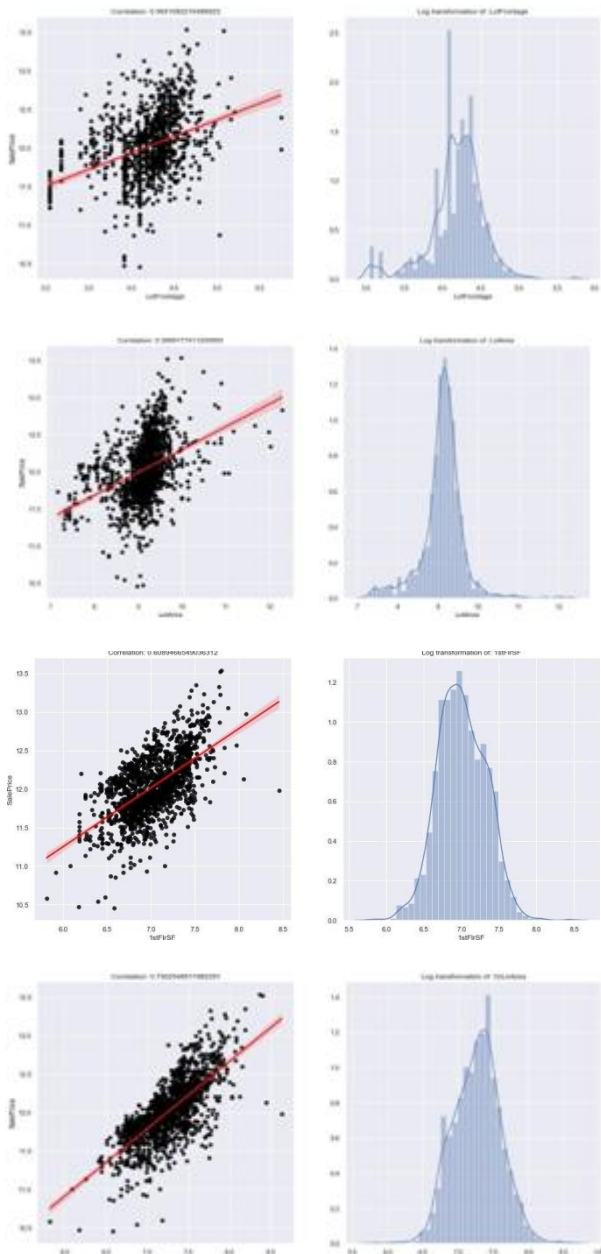


Figure 13: Logarithmic comparison between continuous values and SalePrice

The log transformation revealed interesting insights. While LotArea showed limited correlation with SalePrice,

GrLivArea and 1stFlrSF exhibited strong correlations. However, the presence of zero values prevented log transformation for some features. Further analysis and preprocessing are necessary to address these irregularities and enhance the correlation between the features and the target variable

Categorical Variables

Categorical variables provide insights into the relationship between factors and SalePrice. Analyzing their distribution, impact, and handling of missing values helps identify influential categories and guides feature selection and modeling decisions for accurate predictions.

```

categorical_features = [feature for feature
in train.columns if train[feature].dtypes ==
'O']
print(f"Number of categorical feature: {len
(categorical_features)}")

```

Output:

Number of categorical feature: 43
['MSZoning', 'Street', 'Alley', 'LotShape', 'LandContour', 'Utilities', 'LotConfig', 'LandSlope', 'Neighborhood', 'Condition1', 'Condition2', 'BldgType', 'HouseStyle', 'RoofStyle', 'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2', 'Heating', 'HeatingQC', 'CentralAir', 'Electrical', 'KitchenQual', 'Functional', 'FireplaceQu', 'GarageType', 'GarageFinish', 'GarageQual', 'GarageCond', 'PavedDrive', 'PoolQC', 'Fence', 'MiscFeature', 'SaleType', 'SaleCondition']

To visualize the relationship between Categorical Variables and SalePrice

```

plt.figure(figsize=(40,30))
for i, feature in enumerate(categorical_features, 1):
data = train.copy()
temp = data.groupby(feature)['SalePrice'].median()
plt.subplot(9,5,i)
sns.barplot(temp.index, temp.values)
plt.xticks(rotation=45)
plt.tight_layout(hpad=1.2)
plt.show()

```

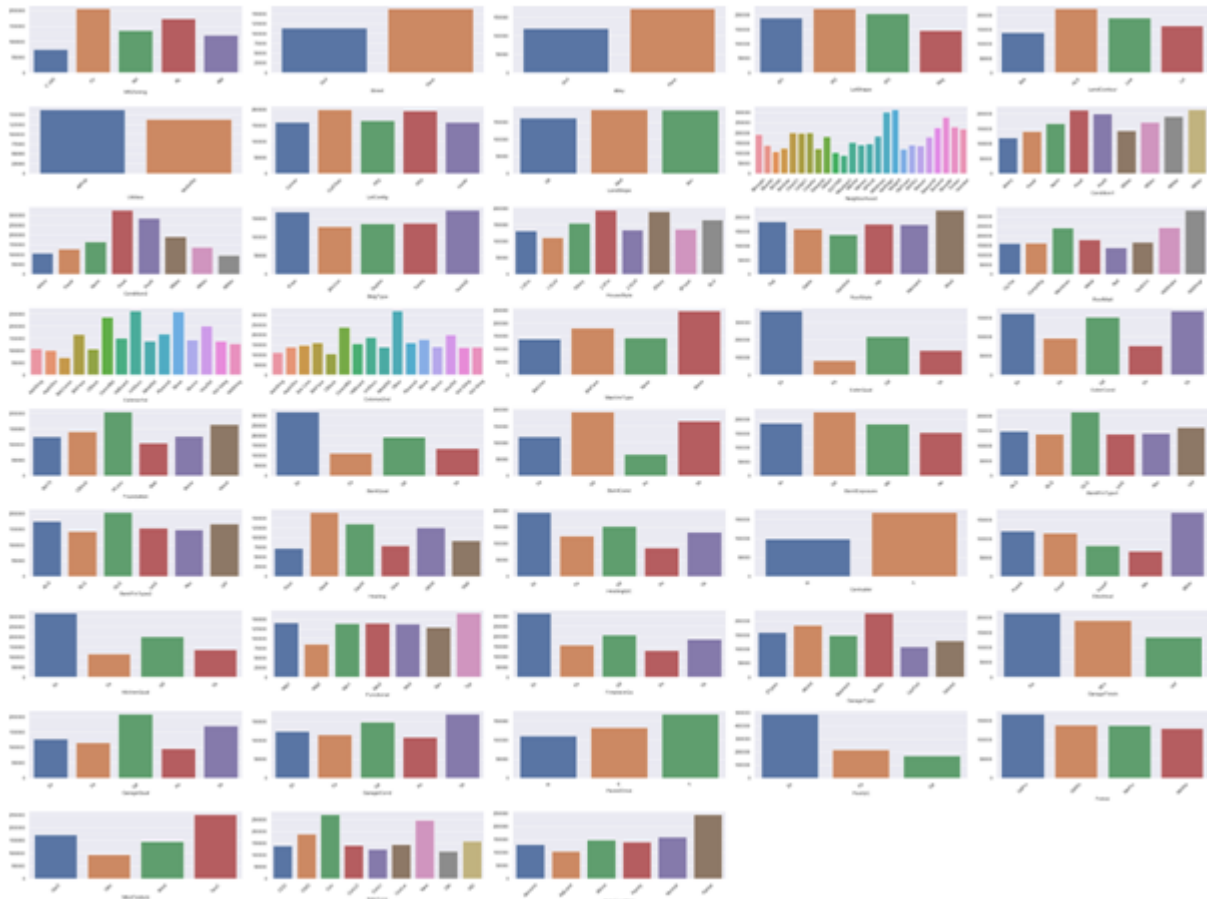


Figure 14: Comparison between Categorical Variables and SalePrice

- Some categorical variables have stronger impacts on the sale price, such as the neighborhood, exterior quality, kitchen quality, and the presence of a fireplace or central air conditioning.
- Other variables may not have significant effects, such as the type of electrical system or certain miscellaneous features.
- For some groups of related variables, selecting a representative feature with high correlation to sale price can be sufficient, such as GarageCars representing garage-related features.
- Overall, careful feature selection and analysis will help determine the most influential categorical variables for accurate modeling.

4. Conclusion

In conclusion, this study delved into the realm of data analysis and feature selection for house price prediction. Through extensive exploration and analysis of the dataset, we uncovered valuable insights regarding the variables that strongly influence the sale price of houses. Key factors such as overall quality, living area, garage capacity, and basement size emerged as significant contributors to the price. By addressing data quality issues like missing values and outliers, and applying feature engineering and transformation techniques, we were able to enhance the interpretability and predictive power of our models.

Looking towards the future, the field of data analysis holds great promise. Advancements in machine learning, artificial

intelligence, and big data processing will pave the way for more sophisticated and accurate predictive models. Integration of diverse data sources, including geospatial and social media data, will further enrich our understanding of the factors shaping house prices. Automated feature selection and model interpretability methods will streamline analysis processes and enable better decision-making.

As we navigate this evolving landscape, it is important to strike a balance between comprehensive analysis and practical implementation. The findings from this study provide a solid foundation for developing robust house price prediction models. By leveraging data-driven insights and advanced analytical techniques, we can make more informed decisions in the real estate industry, empowering stakeholders with accurate pricing information and valuable insights. Ultimately, the potential for transformative advancements in predictive analytics within the real estate domain is substantial.

References

- [1] <https://www.kaggle.com/pmarcelino/comprehensive-data-exploration-with-python>
- [2] <https://www.kaggle.com/code/juliencs/a-study-on-regression-applied-to-the-ames-dataset/notebook>
- [3] <https://www.kaggle.com/code/siddheshpujari/eda-and-prediction-of-house-price/notebook>
- [4] Singh, Archana Sharma, Apoorva Dubey, Gaurav. (2020). Big data analytics predicting real estate prices. International Journal of System Assurance Engineering

and Management.11. 10.1007/s13198 - 020 - 00946 - 3.

- [5] Wu, Chao ye, Xinyue Ren, Fu Wan, You Ning, Pengfei Du, Qingyun. (2016). Spatial and Social Media Data Analytics of Housing Prices in Shenzhen, China. PLOS ONE.11.10.1371/journal. pone.0164553.
- [6] <https://www.kaggle.com/code/apapiu/regularized-linear-models/notebook>
- [7] Li, Mingzhao Bao, Zhifeng Sellis, Timos Yan, Shi Zhang, Rui. (2018). HomeSeeker: A Visual Analytics System of Real Estate Data. Journal of Visual Languages Computing.45.10.1016/j. jvlc.2018.02.001.