

```
with ZipFile(zip_path) as zip_ref:  
    zip_ref.extractall(extraction_path)  
  
# Load the dataset into a Pandas DataFrame  
csv_file_path = os.path.join(extraction_path, 'train.csv')  
train_data = pd.read_csv(csv_file_path)  
  
# Display the first few rows of the DataFrame to ensure it loaded correctly  
print(train_data.head())  
train_data
```

```
Out[2]:
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	...	PoolArea	Price
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	...	0	1455
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	...	0	1455
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	...	0	1455
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	...	0	1455
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	...	0	1455
...
1455	1456	60	RL	62.0	7917	Pave	NaN	Reg	Lvl	AllPub	...	0	1455
1456	1457	20	DJ	85.0	13175	Dave	NaN	Den	Lvl	AllPub	...	0	1455

```

0 1 MSSubClass LotFrontage LotArea OverallQual OverallCond YearBuilt \
1 1 60 65.0 8450 7 5 2003
2 2 20 80.0 9600 6 8 1976
3 3 60 68.0 11250 7 5 2001
3 4 70 60.0 9550 7 5 1915
5 5 60 84.0 14260 8 5 2000

YearRemodAdd MasVnrArea BsmtFinSFl ... GarageType GarageFinish \
0 2003 196.0 706 ... Attchd Rfn
1 1976 0.0 978 ... Attchd Rfn
2 2002 162.0 486 ... Attchd Rfn
3 1970 0.0 216 ... Detchd UnF
4 2000 350.0 655 ... Attchd UnF

GarageQual GarageCond PavedDrive PoolQC Fence MiscFeature SaleType \
0 TA TA Y NaN NaN NaN WO
1 TA TA Y NaN NaN NaN WO
2 TA TA Y NaN NaN NaN WO
3 TA TA Y NaN NaN NaN WO
4 TA TA Y NaN NaN NaN WO

SaleCondition
0 Normal
1 Normal
2 Normal
3 Abnrmal
4 Normal

[5 rows x 81 columns]

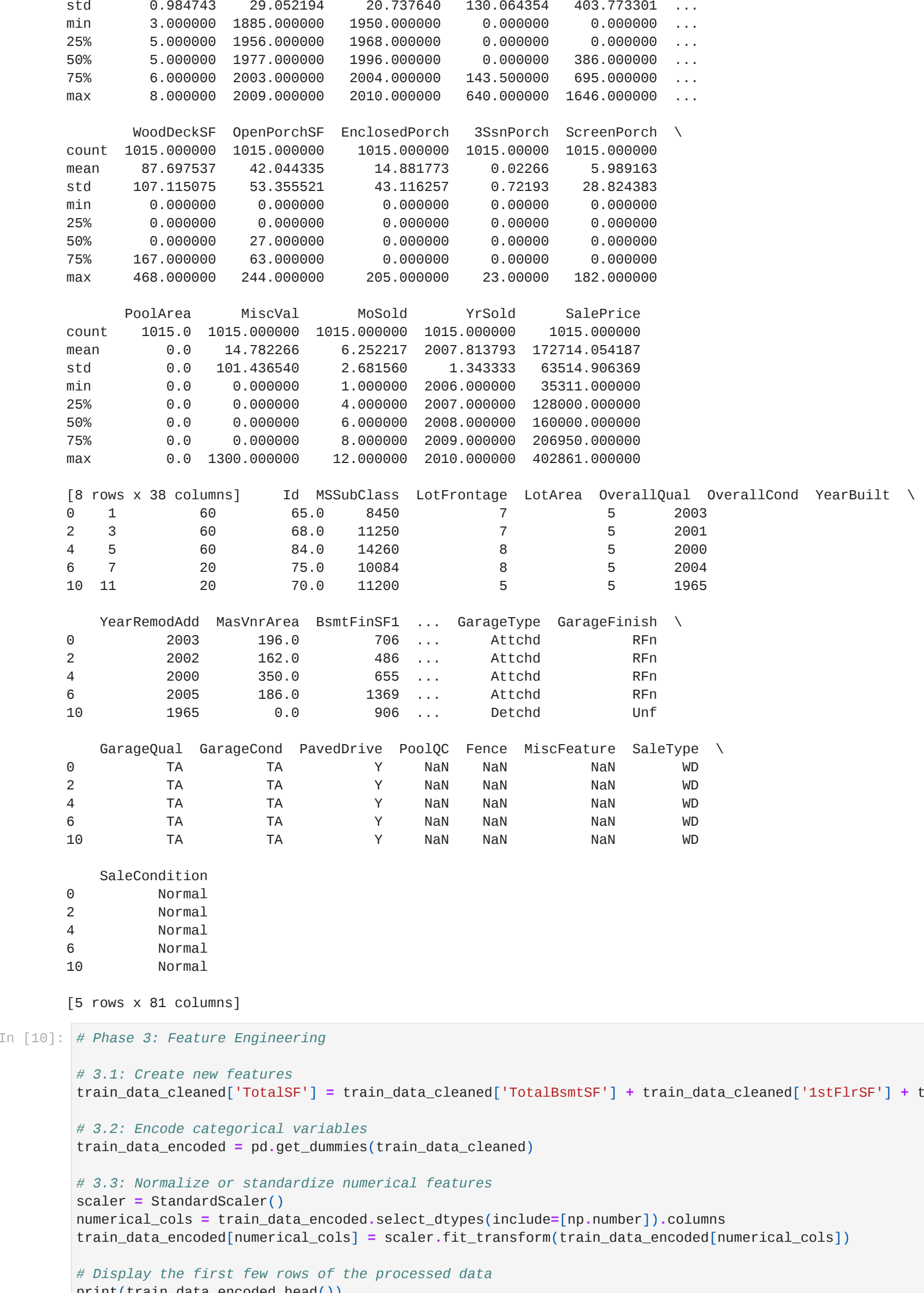
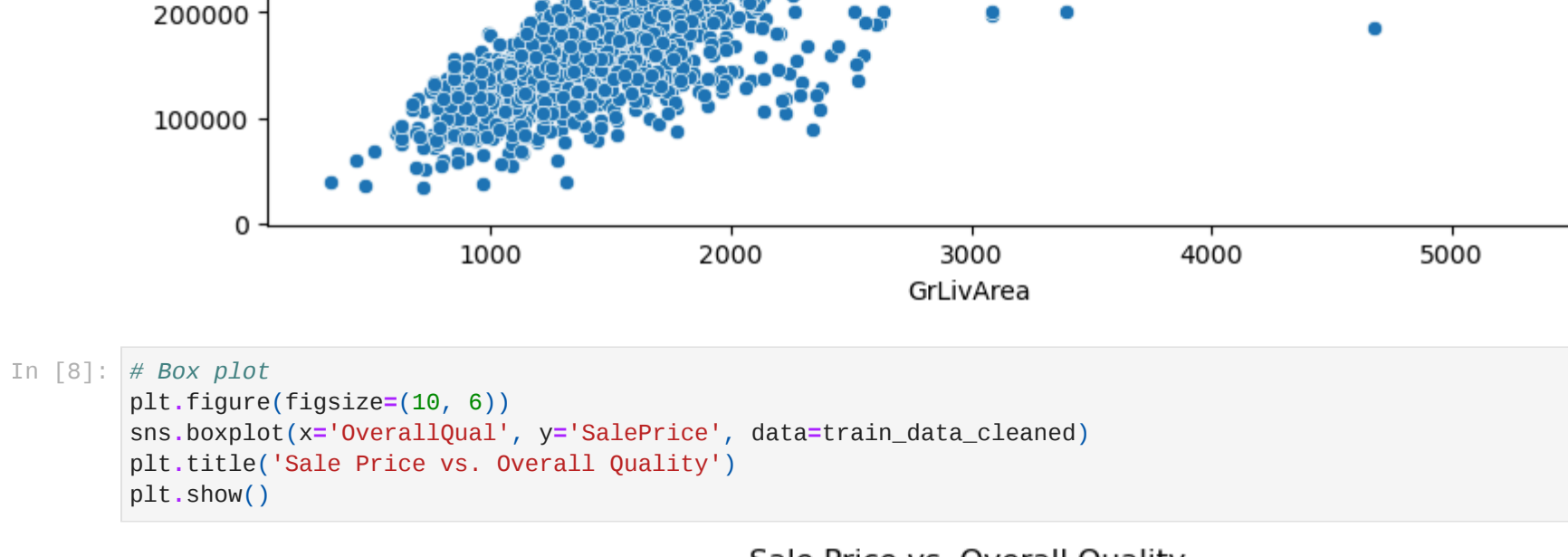
In [5]: # Perform data cleaning to ensure the dataset is ready for analysis
initial_rows = train_data.shape[0]
train_data_cleaned = train_data.fillna(copy())

# Display the first few rows of the cleaned data
print(train_data_cleaned.head(), initial_rows, missing_values.head(10))

0 1 MSSubClass LotFrontage LotArea OverallQual OverallCond YearBuilt \
0 1 60 65.0 8450 7 5 2003
1 2 20 80.0 9600 6 8 1976
2 3 60 68.0 11250 7 5 2001
3 4 70 60.0 9550 7 5 1915
4 5 60 84.0 14260 8 5 2000

```

A histogram showing the distribution of the number of children per family. The x-axis is labeled 'children' and ranges from 0 to 10. The y-axis is labeled 'count' and ranges from 0 to 150. The distribution is roughly bell-shaped, peaking at 3 children with a count of approximately 145.



```

    }

    grid_search = GridSearchCV(RandomForestRegressor(random_state=42), param_grid, cv=3, scoring='neg_mean
    grid_search.fit(X_train, y_train)

    best_rf_model = grid_search.best_estimator_

In [16]: # 4.5: Select and evaluate the best-performing model on the testing set
y_pred_best = best_rf_model.predict(X_test)
rmse_best = mean_squared_error(y_test, y_pred_best, squared=False)

In [17]: # Summary results
rmse_results = {
    "Linear Regression RMSE": rmse_lr,
    "Random Forest RMSE": rmse_rf,
    "Best Model RMSE": rmse_best,
}

```

```
print(rmse_results, best_rf_model)
```

```
{'Linear Regression RMSE': 139260627752.467, 'Random Forest RMSE': 0.29934269764909005, 'Best Model RMSE': 0.29934269764909005} RandomForestRegressor(random_state=42)
```