



# Presentation on Houses Price

Team Project

**Artur**

**Yu Dzhin**

**Shabdan**

**This dataset was used in data science competitions in Kaggle**



**The House Price dataset is a well-known dataset used in machine learning and data analysis. The dataset contains information on residential homes in Ames, Iowa, and includes various features such as the sale price, location, size, age, and other important characteristics of the property. This dataset is often used as a benchmark in machine learning competitions and is frequently used by researchers in the field of real estate analysis.**

**The dataset was compiled by Dean De Cock in 2011 and was intended to be an improvement over similar datasets that were available at the time. It includes 80 explanatory variables that describe various aspects of residential homes in Ames, Iowa, and a target variable, which is the sale price of each home. The dataset is relatively large, containing 1,460 observations, and is commonly used for regression and classification tasks.**

**The House Price dataset is a valuable resource for machine learning practitioners, data analysts, and researchers interested in real estate analysis. Its widespread use and popularity make it an ideal starting point for exploring the application of machine learning techniques to real-world problems.**

Dataset	Step 1	Step 2	Step 3	Step 4			Step 5	Step 6
	Data analysis	Data preprocessing & Feature engineering	Feature Selection	Pipeline		Metric score	Conclusion	
HousePrice CSV	HousePrice CSV	HousePrice CSV	HousePrice CSV	Stacks of techniques	Model	Determine best parameters	HousePrice CSV	HousePrice CSV
				<ul style="list-style-type: none"><li>CategoricalImputer<ul style="list-style-type: none"><li>frequent</li><li>missing</li></ul></li><li>MeanMedianImputer</li><li>RandomSampleImputer</li><li>RareLabelEncoder</li><li>OneHotEncoder</li><li>MeanEncoder</li><li>EqualFrequencyDiscretiser</li><li>DropConstantFeatures</li><li>DropDuplicateFeatures</li><li>SmartCorrelatedSelection</li></ul>	Decision tree classifier	GridSearchCV		
				<ul style="list-style-type: none"><li>CategoricalImputer<ul style="list-style-type: none"><li>frequent</li><li>missing</li></ul></li><li>MeanMedianImputer</li><li>RandomSampleImputer</li><li>RareLabelEncoder</li><li>OneHotEncoder</li><li>MeanEncoder</li><li>EqualFrequencyDiscretiser</li><li>DropConstantFeatures</li><li>DropDuplicateFeatures</li><li>SmartCorrelatedSelection</li></ul>	Random forest classifier	GridSearchCV		
				<ul style="list-style-type: none"><li>CategoricalImputer<ul style="list-style-type: none"><li>frequent</li><li>missing</li></ul></li><li>MeanMedianImputer</li><li>RandomSampleImputer</li><li>RareLabelEncoder</li><li>OneHotEncoder</li><li>MeanEncoder</li><li>EqualFrequencyDiscretiser</li><li>DropConstantFeatures</li><li>DropDuplicateFeatures</li><li>SmartCorrelatedSelection</li></ul>	Support vector classifier	GridSearchCV		
				<ul style="list-style-type: none"><li>CategoricalImputer<ul style="list-style-type: none"><li>frequent</li><li>missing</li></ul></li><li>MeanMedianImputer</li><li>RandomSampleImputer</li><li>RareLabelEncoder</li><li>OneHotEncoder</li><li>MeanEncoder</li><li>EqualFrequencyDiscretiser</li><li>DropConstantFeatures</li><li>DropDuplicateFeatures</li><li>SmartCorrelatedSelection</li></ul>	Logistic regression	GridSearchCV		
				<ul style="list-style-type: none"><li>CategoricalImputer<ul style="list-style-type: none"><li>frequent</li><li>missing</li></ul></li><li>MeanMedianImputer</li><li>RandomSampleImputer</li><li>RareLabelEncoder</li><li>OneHotEncoder</li><li>MeanEncoder</li><li>EqualFrequencyDiscretiser</li><li>DropConstantFeatures</li><li>DropDuplicateFeatures</li><li>SmartCorrelatedSelection</li></ul>	XGBoost Classifier	GridSearchCV		
				<ul style="list-style-type: none"><li>CategoricalImputer<ul style="list-style-type: none"><li>frequent</li><li>missing</li></ul></li><li>MeanMedianImputer</li><li>RandomSampleImputer</li><li>RareLabelEncoder</li><li>OneHotEncoder</li><li>MeanEncoder</li><li>EqualFrequencyDiscretiser</li><li>DropConstantFeatures</li><li>DropDuplicateFeatures</li><li>SmartCorrelatedSelection</li></ul>	AdaBoost Classifier	GridSearchCV		
				<ul style="list-style-type: none"><li>CategoricalImputer<ul style="list-style-type: none"><li>frequent</li><li>missing</li></ul></li><li>MeanMedianImputer</li><li>RandomSampleImputer</li><li>RareLabelEncoder</li><li>OneHotEncoder</li><li>MeanEncoder</li><li>EqualFrequencyDiscretiser</li><li>DropConstantFeatures</li><li>DropDuplicateFeatures</li><li>SmartCorrelatedSelection</li></ul>				

Import Basic Libraries:	Import Imputation Techniques	Import Encoding Techniques
<b>Pandas ,</b> <b>numpy,</b> <b>matplotlib.pyplot,</b> <b>seaborn,</b> <b>train_test_split</b>	<b>Random Sample Imputer</b> <b>Mean Median Imputer</b> <b>Arbitrary Number Imputer</b> <b>Categorical Imputer</b>	<b>PRatioEncoder</b> <b>OneHotEncoder</b> <b>WoEncoder</b> <b>MeanEncoder</b> <b>OrdinalEncoder</b> <b>RareLabelEncoder</b>
<b>Discretization Techniques</b>  <b>EqualWidthDiscretization</b>  <b>EqualFrequencyDiscretizer</b>	<b>Classifiers</b>  <b>Decision Tree Classifier</b> <b>Random Forest Classifier</b> <b>Ada Boost Classifier</b>  <b>SVC</b>  <b>Logistic Regression</b> <b>XGBClassifier</b>	<b>Metrics</b>  <b>roc_auc_score</b> <b>confusion_matrix,</b> <b>PrecisionRecallDisplay</b>  <b>accuracy_score</b> <b>classification Report</b>

# Compose Element

**GridSearchCV**

**cross\_val\_score**

**PipeLine**

**Feature Magnitude**

**Standard Scaler**

**Normalizer**

**RobustScaler**

**Feature Selection**

**mutual\_info\_classif**

**mutual\_info\_regression**

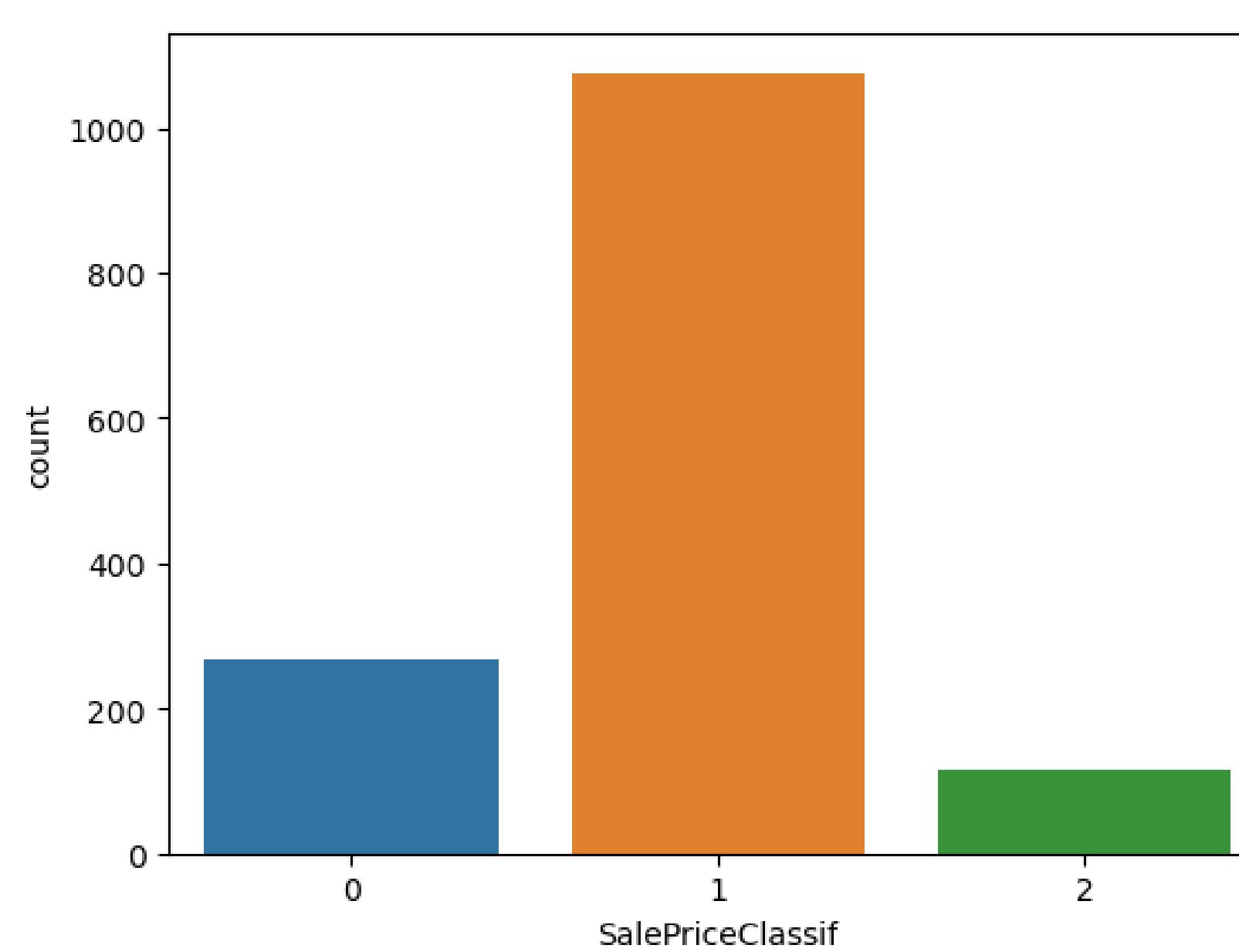
**SelectKBest**

**SelectPercentile**

**DropDuplicatedFeatures**

**DropConstantFeatures**

**SmartCorrealtedSelection**



# Make Classification

```
1 for var in data.columns:  
2     if data[var].isnull().sum() > 0:  
3         print(var, ' Percentage of missing labels: {:.3f} %'.format(data[var].isnull().mean() * 100))
```

LotFrontage Percentage of missing labels: 17.740 %

Alley Percentage of missing labels: 93.767 %

MasVnrType Percentage of missing labels: 0.548 %

MasVnrArea Percentage of missing labels: 0.548 %

BsmtQual Percentage of missing labels: 2.534 %

BsmtCond Percentage of missing labels: 2.534 %

BsmtExposure Percentage of missing labels: 2.603 %

BsmtFinType1 Percentage of missing labels: 2.534 %

BsmtFinType2 Percentage of missing labels: 2.603 %

Electrical Percentage of missing labels: 0.068 %

FireplaceQu Percentage of missing labels: 47.260 %

GarageType Percentage of missing labels: 5.548 %

GarageYrBlt Percentage of missing labels: 5.548 %

GarageFinish Percentage of missing labels: 5.548 %

GarageQual Percentage of missing labels: 5.548 %

GarageCond Percentage of missing labels: 5.548 %

PoolQC Percentage of missing labels: 99.521 %

Fence Percentage of missing labels: 80.753 %

# Checking Missing Values

```
1 # Determine variables with high percentage of missing observations
2 drop_features_missing = [
3     var for var in data.columns if data[var].isnull().mean() > 0.5 and var != 'SalePrice'
4 ]
5
6 print('Count of missing values \n')
7 for var in drop_features_missing:
8     print(var, ': ', data[var].isnull().sum())
```

Count of missing values

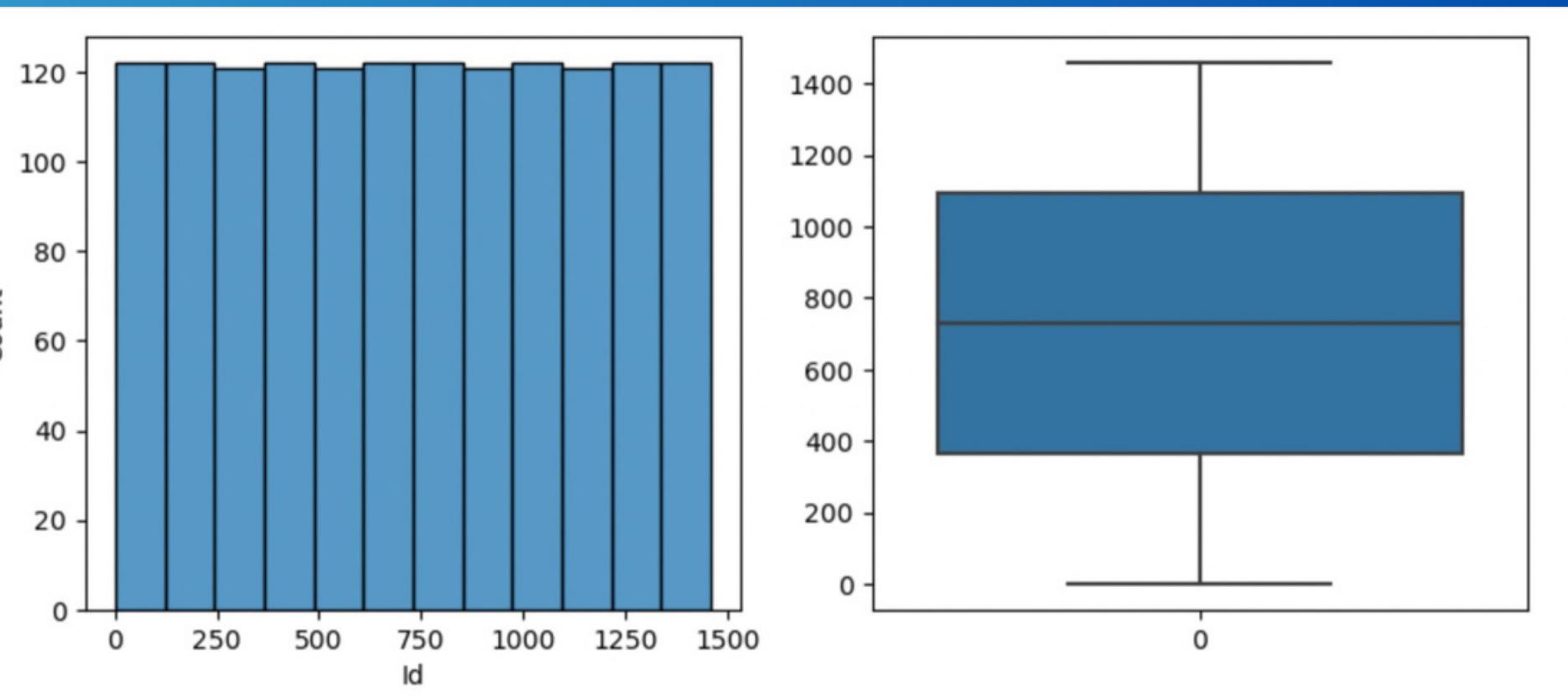
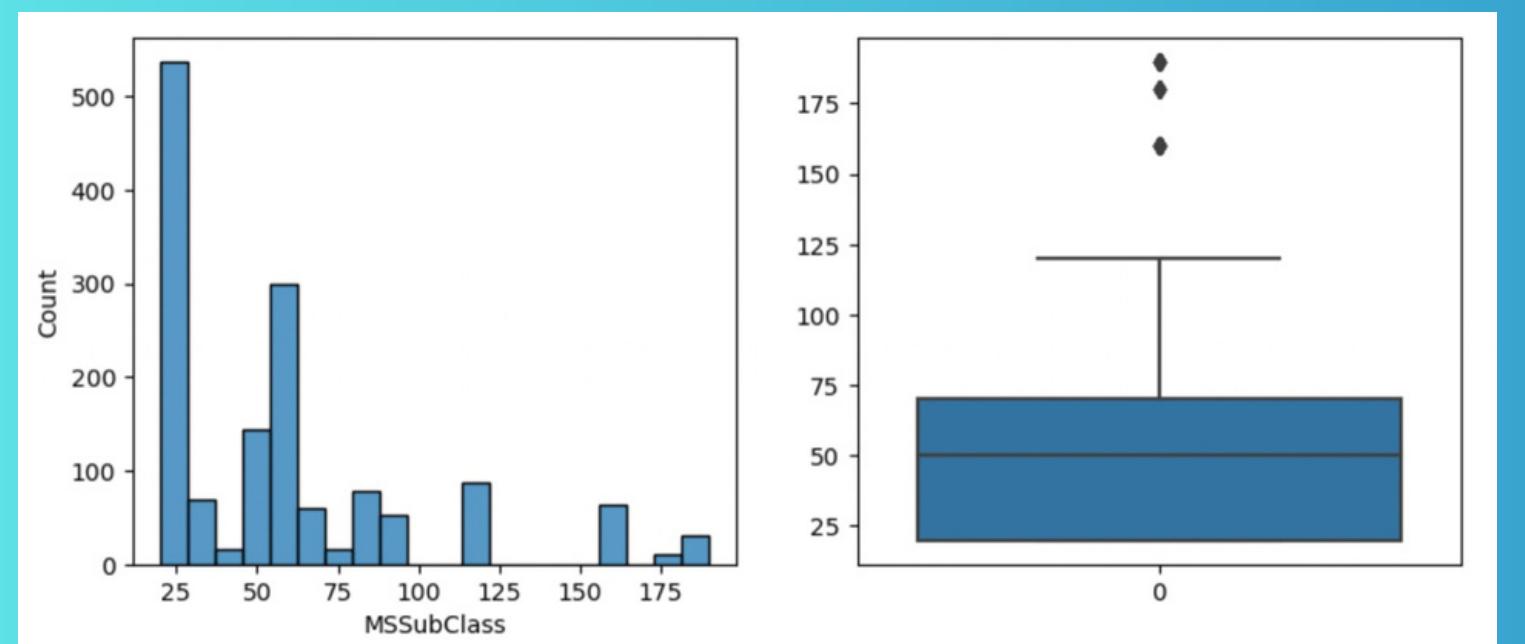
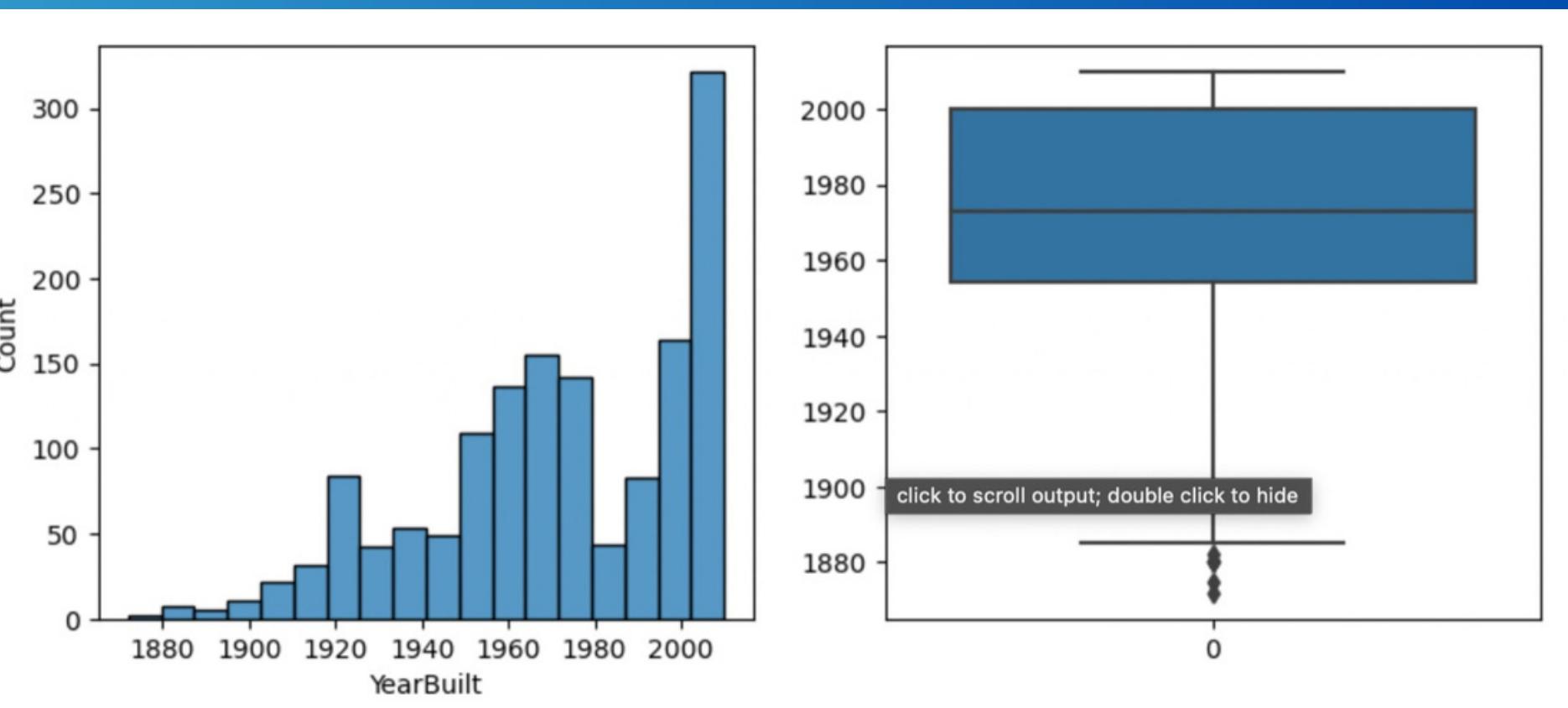
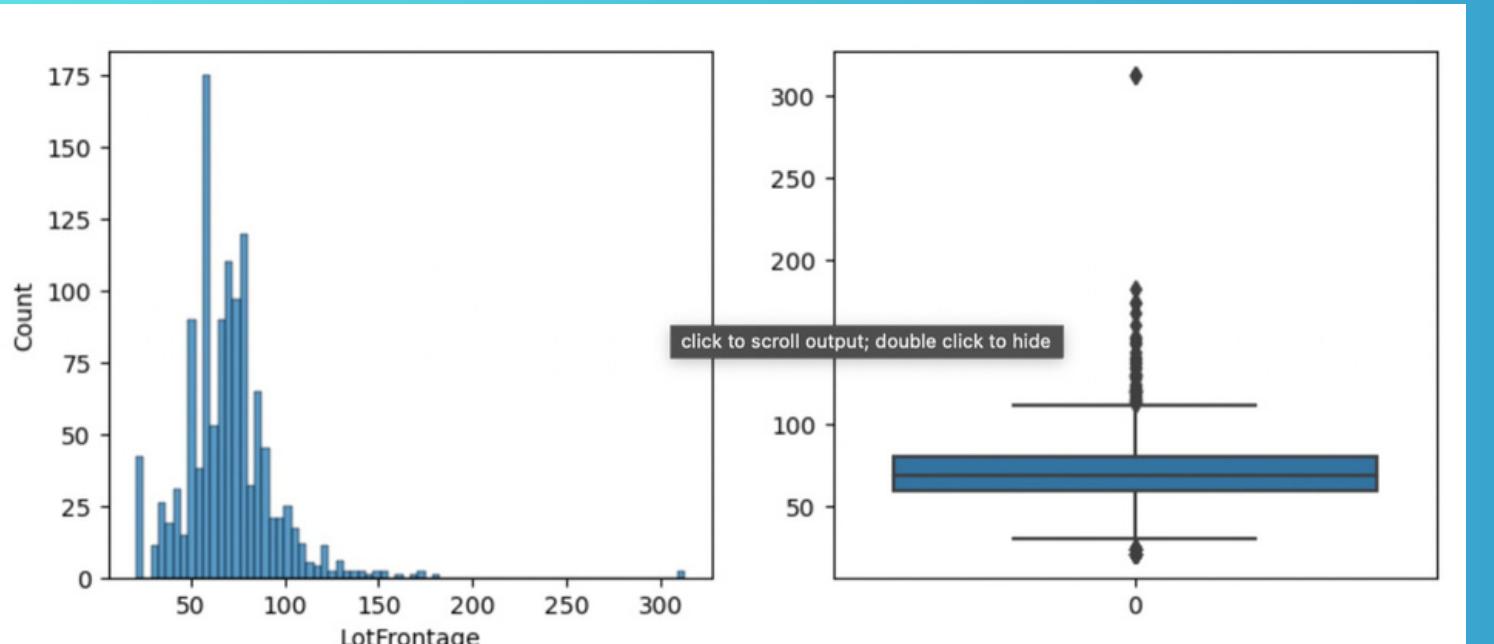
```
Alley : 1369
PoolQC : 1453
Fence : 1179
MiscFeature : 1406
```

## Drop variables with a lot of missing observations

**variables with a lot of missing  
observations do not give us enough information**

## 3.4 Distribution and outliers

```
1 for var in continuous_features_data:  
2     sns.histplot(data[var])  
3     plt.show()
```



## 4. Feature-engineering

```
1 # Splitting dataset on train and test sets
2 X_train, X_test, y_train, y_test = train_test_split(
3     data.drop(labels = 'SalePriceClassif', axis = 1),
4     data['SalePriceClassif'],
5     test_size = 0.3
6 )
7
8 X_train.shape, X_test.shape
```

```
((1022, 75), (438, 75))
```

```
1 X_train.head(10)
```

	MSSubClass	MSZoning	LotFrontage	LotArea	Street	LotShape	LandContour	Utilities	LotConfig	LandSlope	Neighborhood	Condition1	Condition2	E
453	60	FV	75.0	9000	Pave	Reg	Low	AllPub	Inside	Gtl	Somerst	Norm	Norm	Ex
987	20	RL	83.0	10159	Pave	IR1	Low	AllPub	Inside	Gtl	NridgHt	Norm	Norm	Ex
715	20	RL	78.0	10140	Pave	Reg	Low	AllPub	Inside	Gtl	NWAmes	Norm	Norm	Ex
1041	60	RL	NaN	9130	Pave	Reg	Low	AllPub	Inside	Gtl	NWAmes	Feedr	Norm	Ex
402	30	RL	60.0	10200	Pave	Reg	Low	AllPub	Inside	Gtl	Sawyer	Norm	Norm	Ex
741	20	RL	65.0	6768	Pave	IR1	Low	AllPub	Inside	Gtl	Sawyer	Feedr	Norm	Ex
246	190	RM	69.0	9142	Pave	Reg	Low	AllPub	Inside	Gtl	OldTown	Norm	Norm	Ex
1230	90	RL	NaN	18890	Pave	IR1	Low	AllPub	Inside	Gtl	Sawyer	Feedr	RR Ae	Ex
229	120	RL	43.0	3182	Pave	Reg	Low	AllPub	Inside	Gtl	Blmgtn	Norm	Norm	Ex
1448	50	RL	70.0	11767	Pave	Reg	Low	AllPub	Inside	Gtl	Edwards	Norm	Norm	Ex

# Feature Engineering and Data Separation to X\_train and X\_test sets

## 4.1 Create new features

```
1 X_train = dateTimeHandling(X_train)
2 X_test = dateTimeHandling(X_test)
3
4 X_train.shape, X_test.shape
```

```
((1022, 74), (438, 74))
```

```
1 categorical_features_train = [
2     var for var in X_train.columns if X_train[var].dtype == 'O'
3 ]
4
5 numerical_features_train = [
6     var for var in X_train.columns if X_train[var].dtype != 'O' and var not in categorical_features_train
7 ]
8
9 len(numerical_features_train), len(categorical_features_train)
```

```
(35, 39)
```

**Create the new variables (function "dateTimeHandling")**

**Count of Categorical and Numerical features**

## 4.2 Missing data imputation

```
1 need_impute = [  
2     var for var in X_train.columns if X_train[var].isnull().mean() > 0.0  
3 ]  
4  
5 print('Count of features with missing labels: ', len(need_impute))
```

Count of features with missing labels: 15

```
1 for var in need_impute:  
2     print(var)  
3     print('Percentage of missing variables: {:.3f}%'.format(X_train[var].isnull().sum() / len(data) * 100))  
4     print()
```

LotFrontage

Percentage of missing variables: 12.329%

MasVnrType

Percentage of missing variables: 0.274%

MasVnrArea

Percentage of missing variables: 0.274%

BsmtQual

Percentage of missing variables: 1.575%

BsmtCond

Percentage of missing variables: 1.575%

BsmtExposure

Percentage of missing variables: 1.575%

BsmtFinType1

Percentage of missing variables: 1.575%

BsmtFinType2

Percentage of missing variables: 1.644%

Electrical

Percentage of missing variables: 0.068%

FireplaceQu

Percentage of missing variables: 33.767%

GarageType

Percentage of missing variables: 3.699%

GarageYrBlt

Percentage of missing variables: 3.699%

GarageFinish

Percentage of missing variables: 3.699%

GarageQual

Percentage of missing variables: 3.699%

GarageCond

Percentage of missing variables: 3.699%

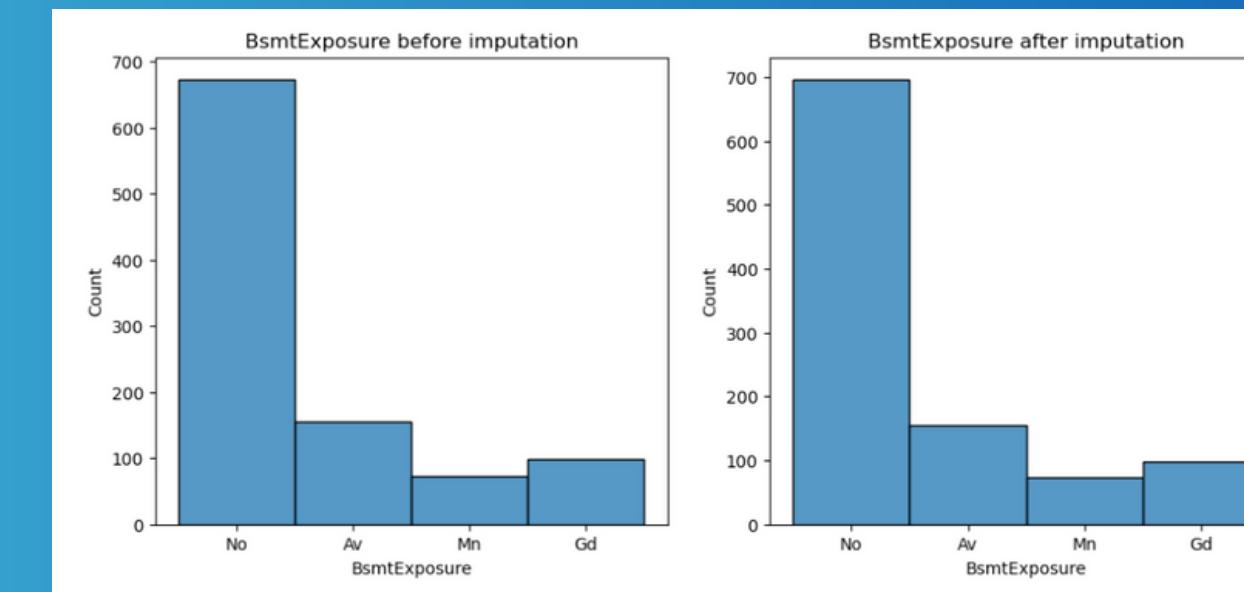
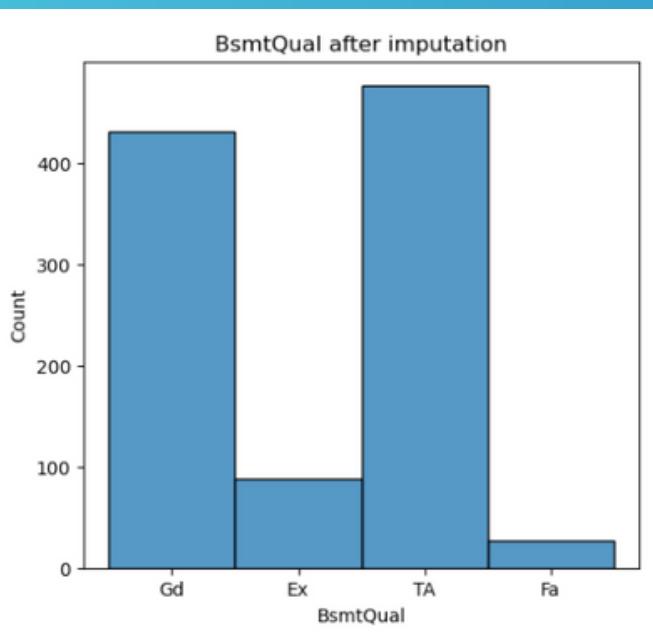
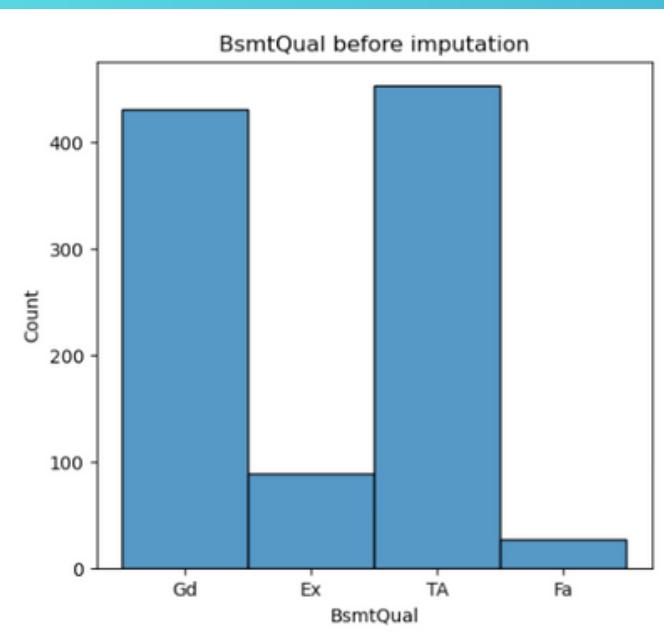
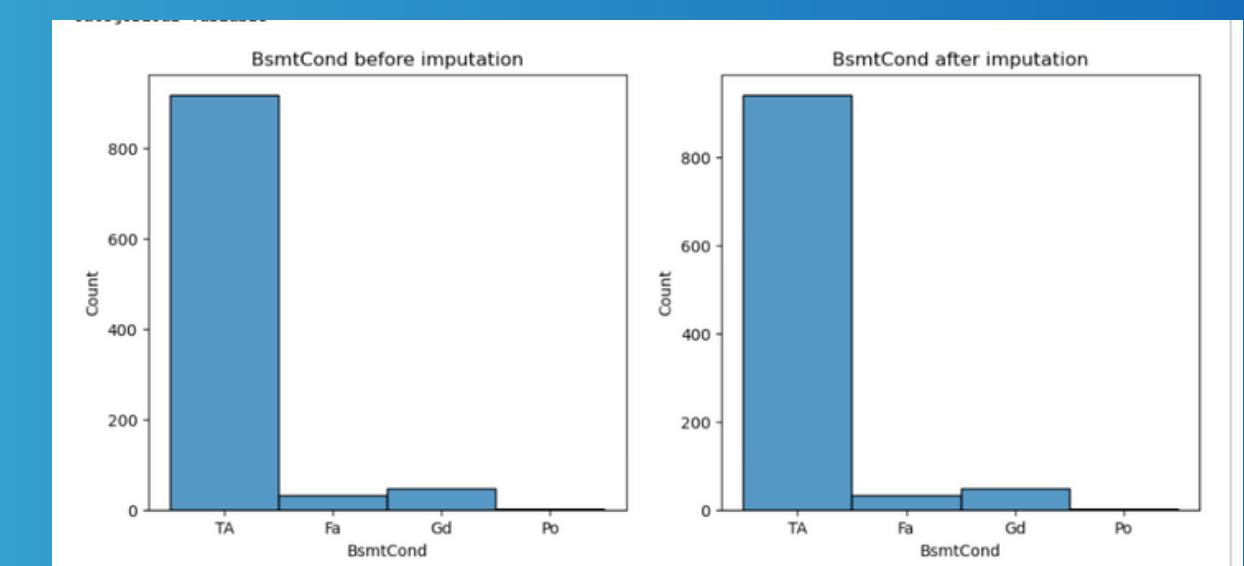
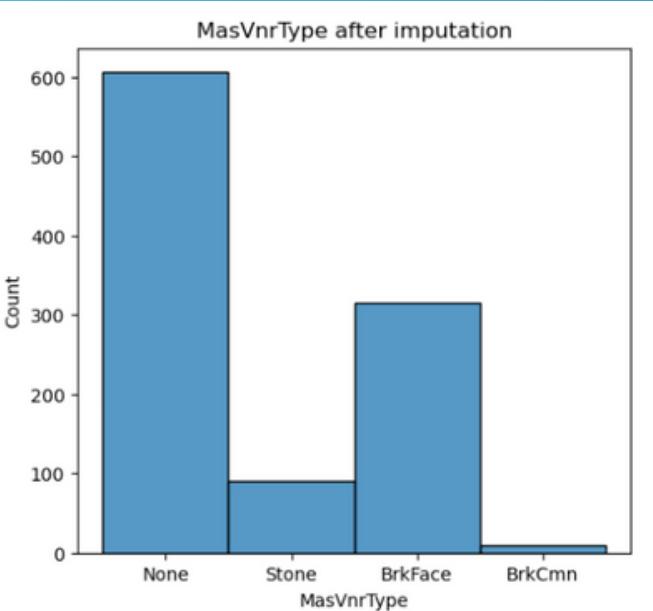
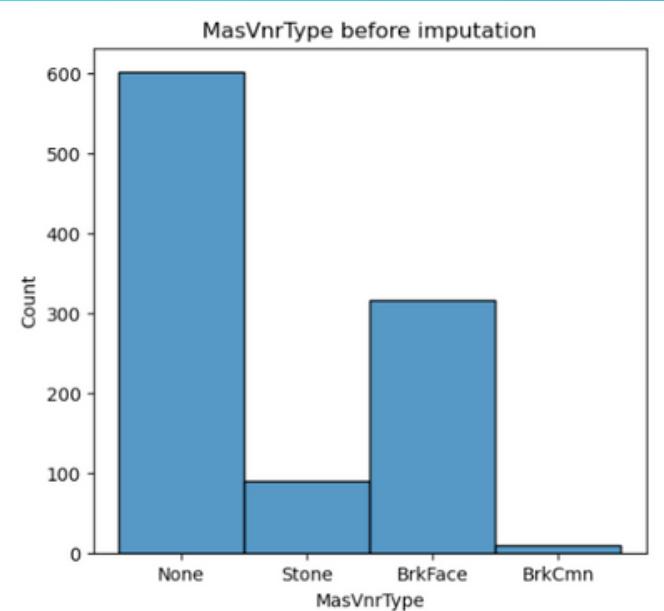
## Categorical imputation methods

### Frequent imputation

```
1 imputer_frequency = CategoricalImputer(imputation_method = 'frequent', variables = need_impute_cat)
2 imputer_frequency.fit(X_train)
3
4 X_train_impute_frequent = imputer_frequency.transform(X_train)
5
6 X_test_impute_frequent = imputer_frequency.transform(X_test)
```

```
1 # Distribution in train set after frequent imputation
2 for var in need_impute_cat:
3     check_dist_and_variance(X_train, X_train_impute_frequent, var)
```

# Frequent Imputation Method

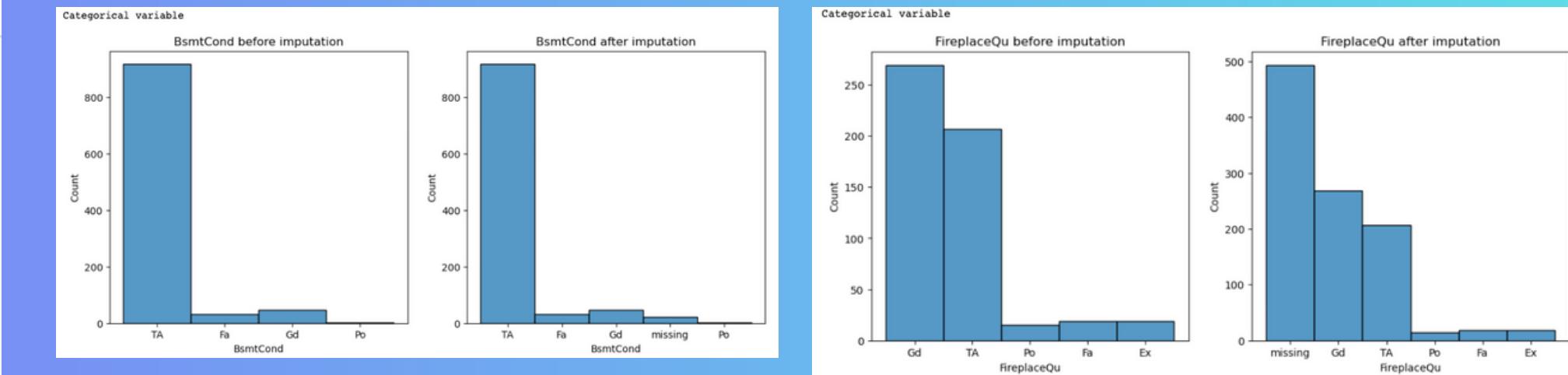


## Missing labels imputation

```

1 imputer_miss_labels = CategoricalImputer(
2     imputation_method = 'missing',
3     fill_value = 'missing',
4     variables = need_impute_cat
5 )
6
7 imputer_miss_labels.fit(X_train)
8
9 X_train_impute_miss_labels = imputer_miss_labels.transform(X_train)
10
11 X_test_impute_miss_labels = imputer_miss_labels.transform(X_test)

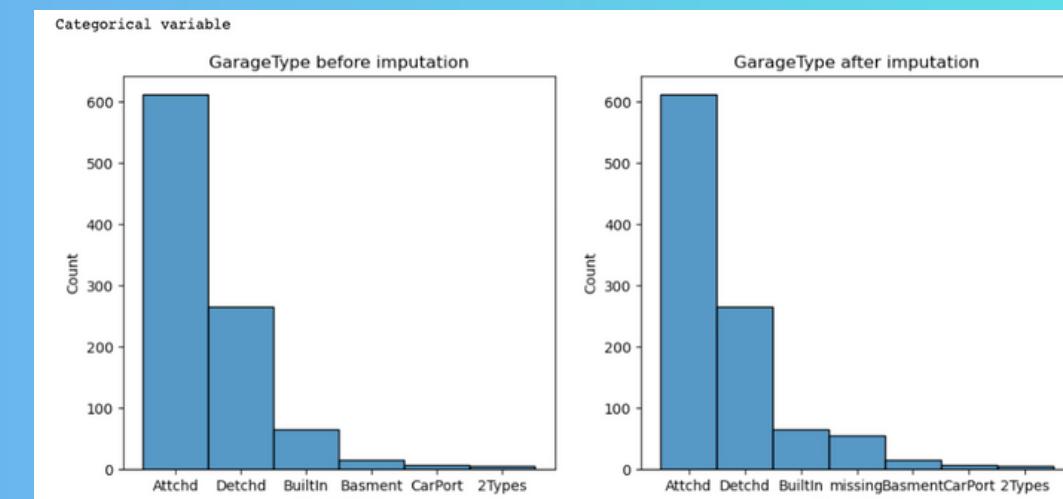
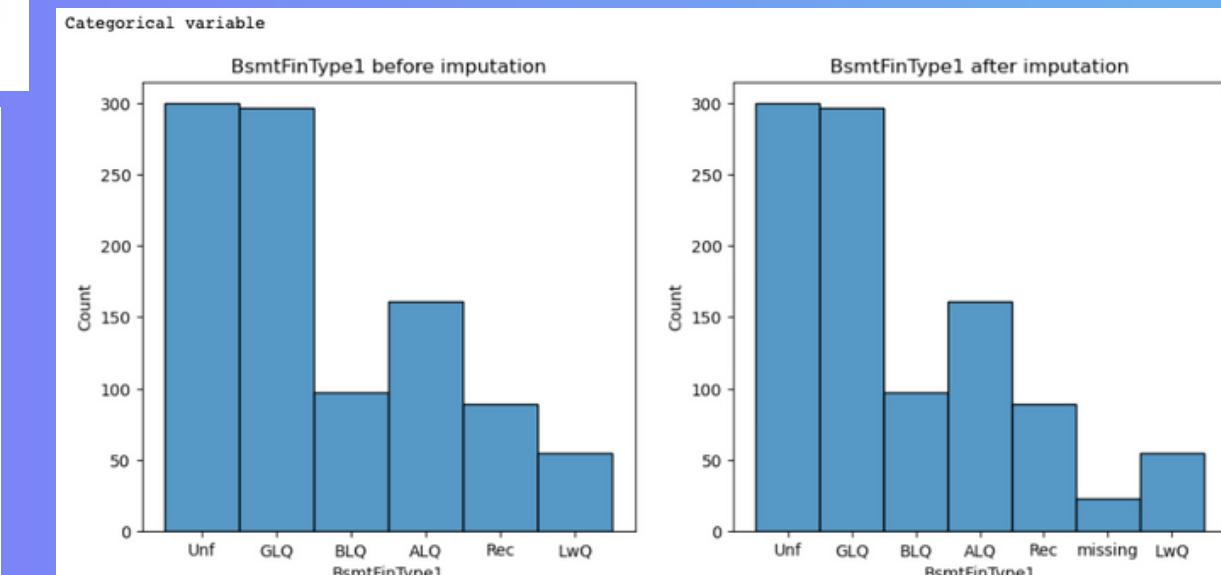
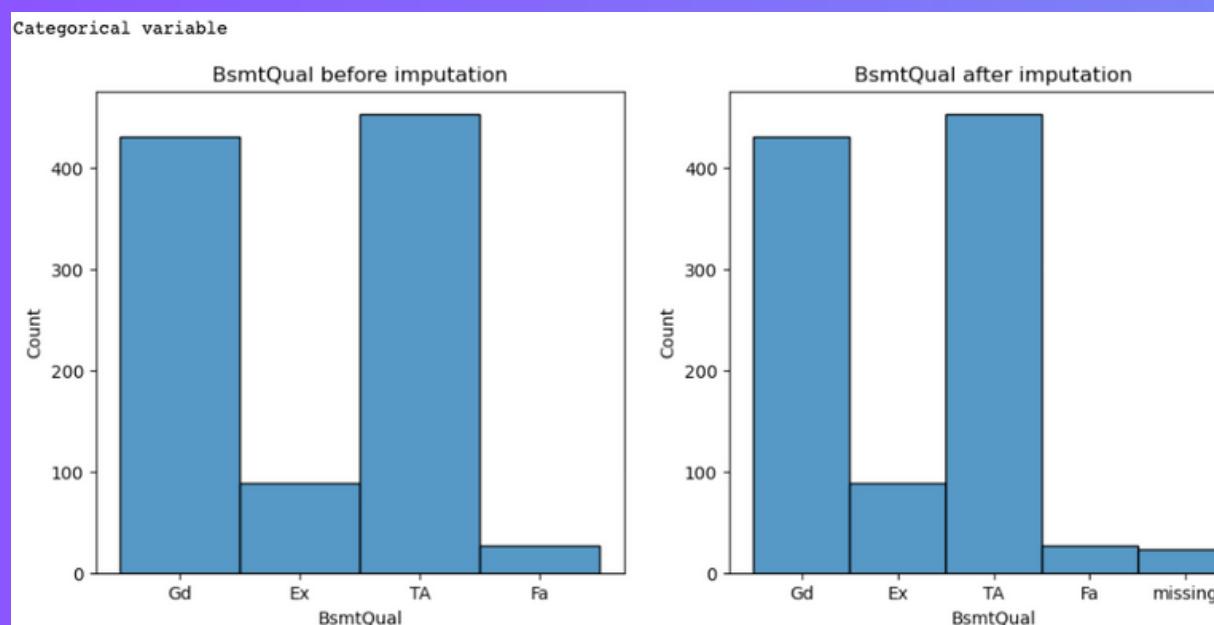
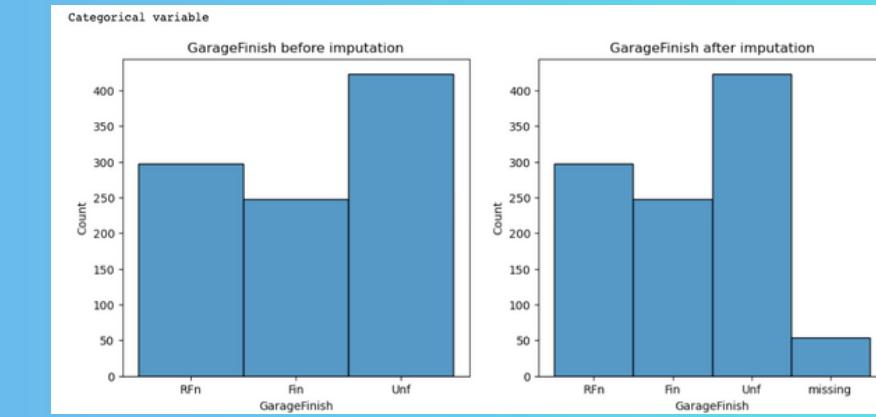
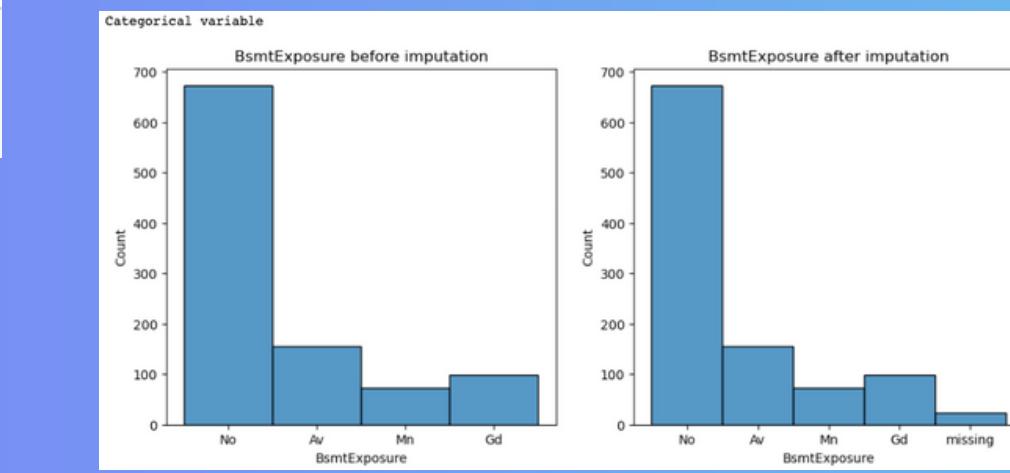
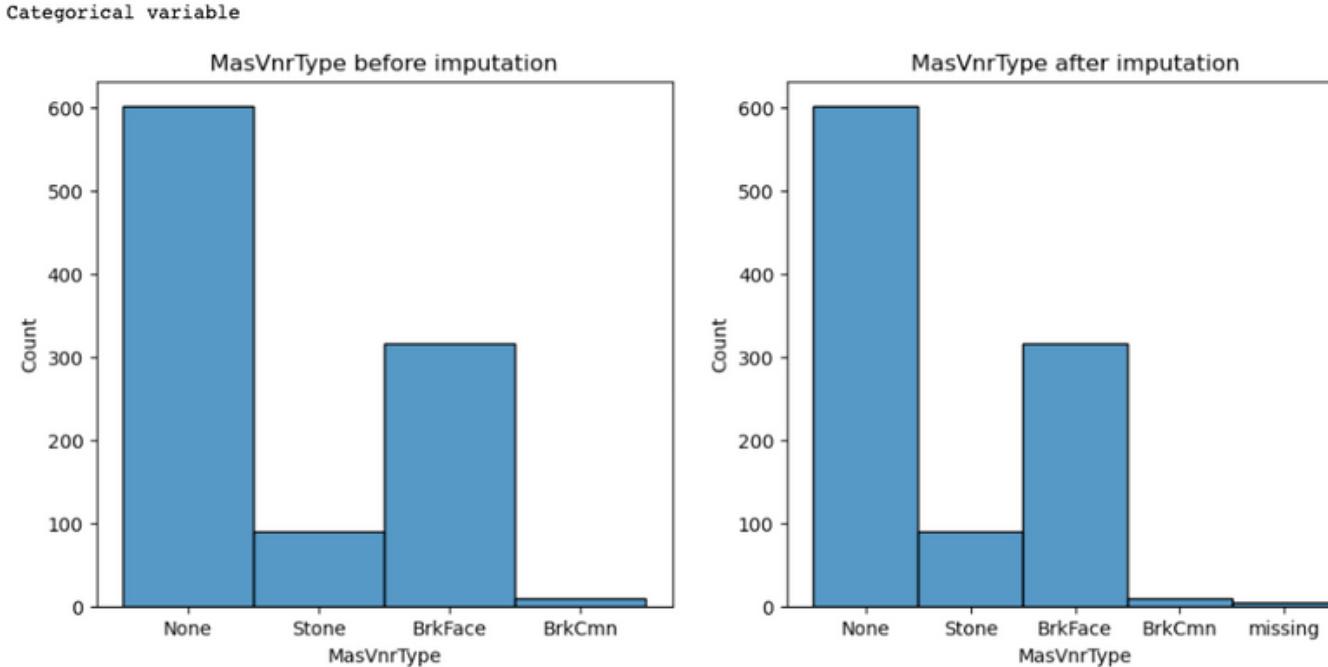
```



```

1 # Check distribution in train set after missing_labels imputation
2 for var in need_impute_cat:
3     check_dist_and_variance(X_train, X_train_impute_miss_labels, var)

```



# Missing Label Imputation

## Numerical imputation methods

### Arbitrary number imputer (arbitrary number = -1)

```
1 imputer_arbitrary = ArbitraryNumberImputer(variables = need_impute_num, arbitrary_number = -1)
2
3 imputer_arbitrary.fit(X_train, y_train)
4
5 X_train_impute_arb_minus = imputer_arbitrary.transform(X_train)
6
7 X_test_impute_arb_minus = imputer_arbitrary.transform(X_test)
```

```
1 # Check distribution after arbitrary number imputation (arbitrary value = -1)
2 for var in need_impute_num:
3
4     check_dist_and_variance(X_train, X_train_impute_arb_minus, var)
```

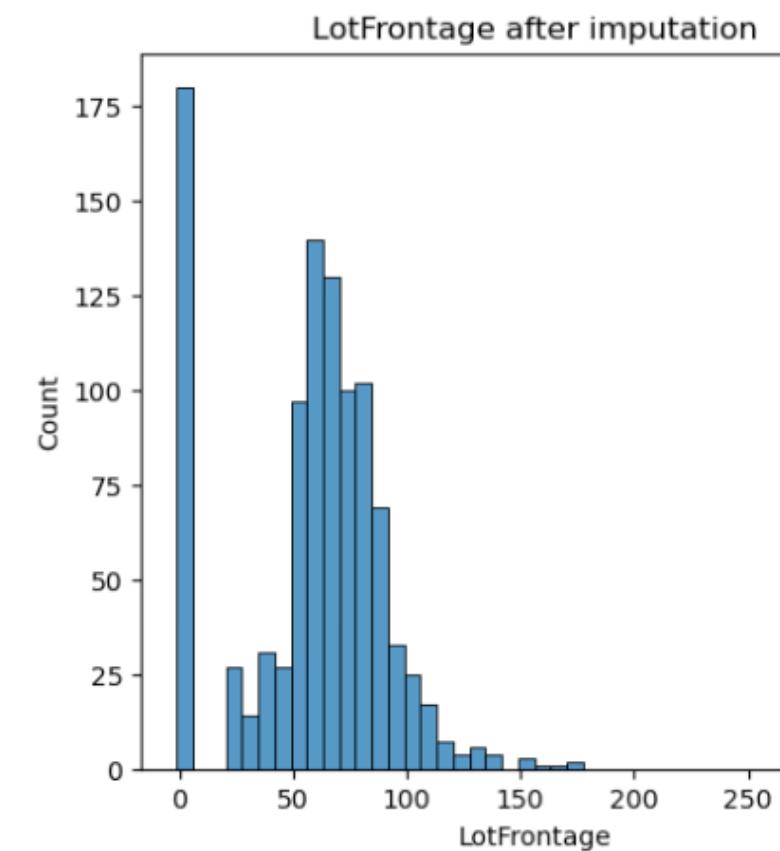
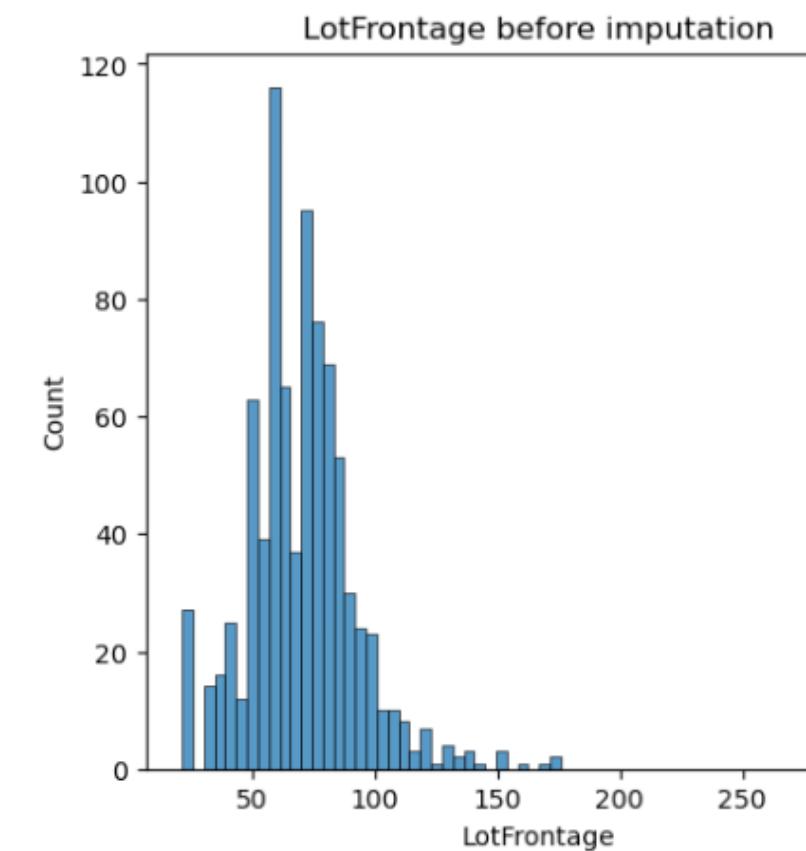
## Numerical imputation Methods

### Arbitrary Number Imputer (arbitrary number = -1)

### Check Distribution after arbitrary number imputation (arbitrary value = -1)

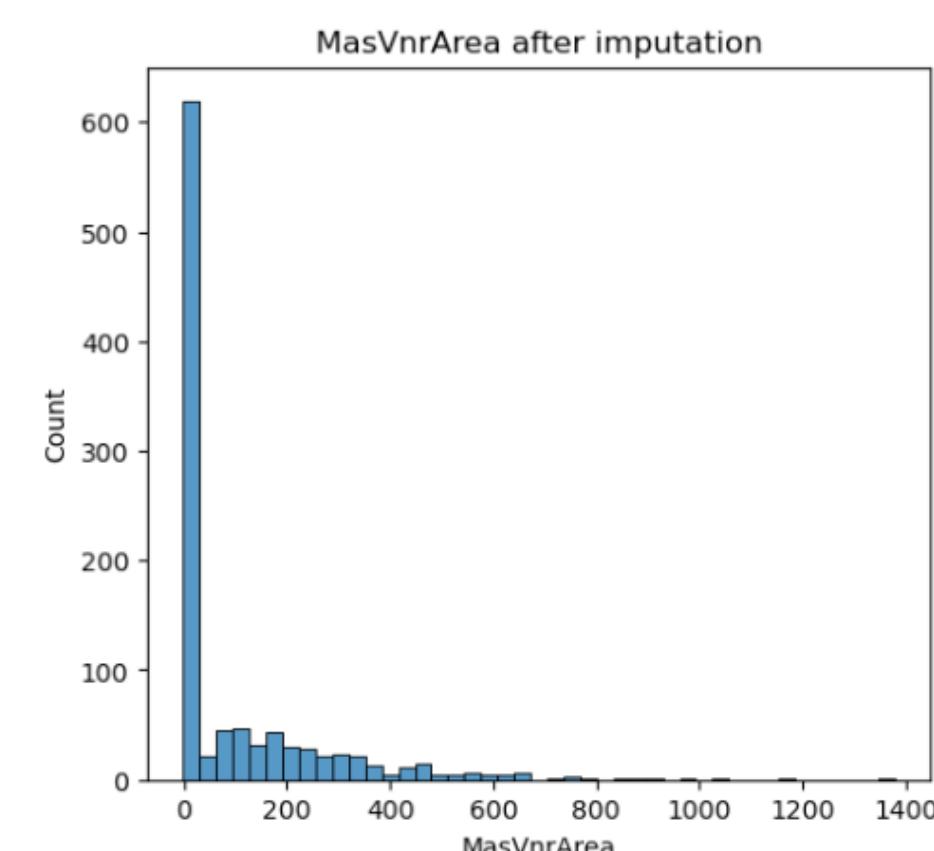
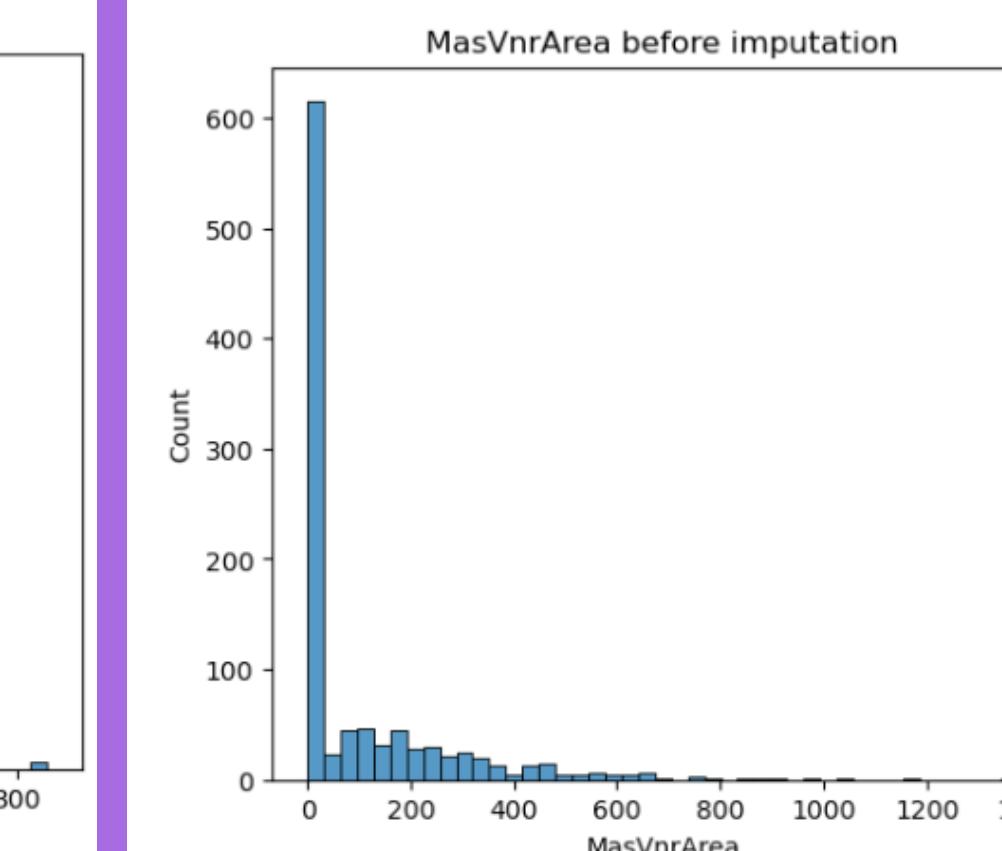
Variance before imputation: 631.2577620805455

Variance after imputation: 1258.8131901305462



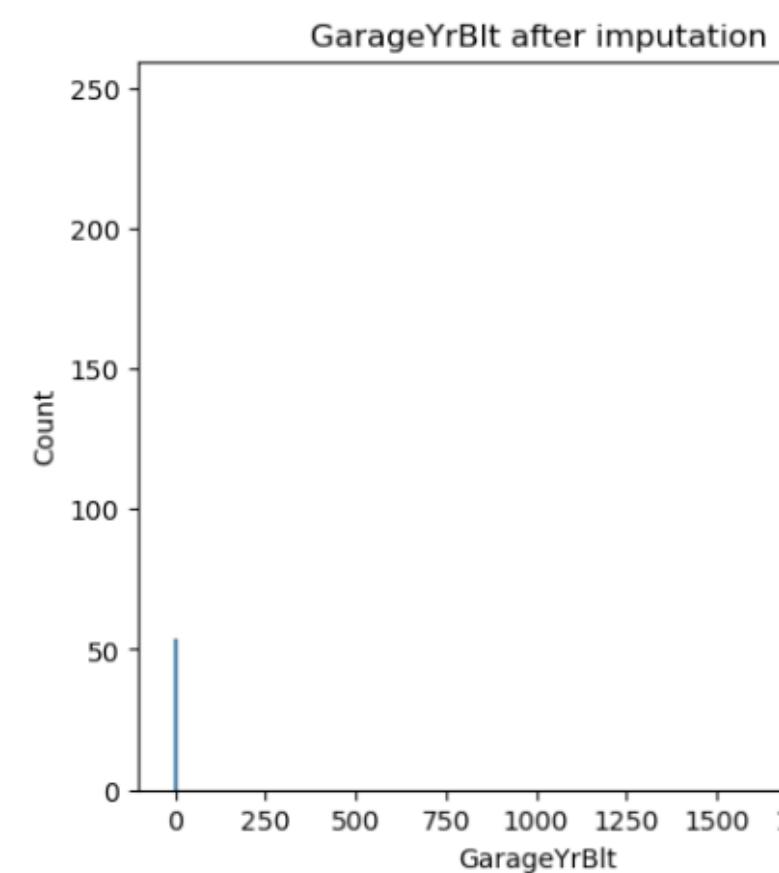
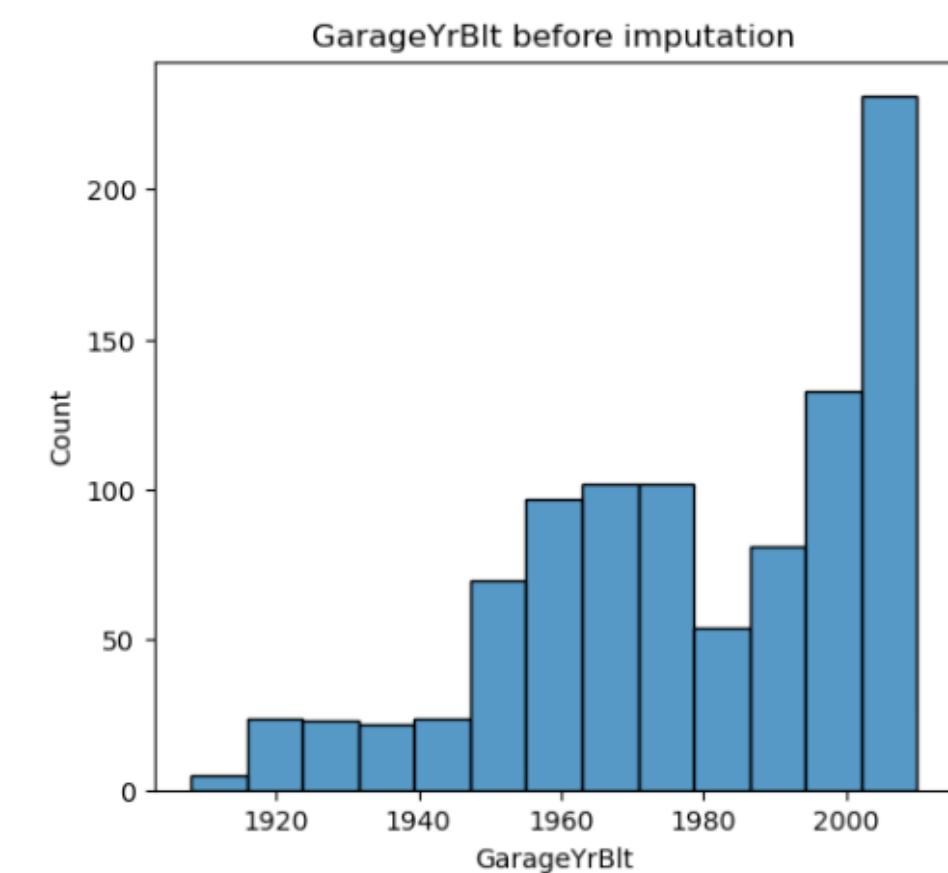
Variance before imputation: 31981.573954946656

Variance after imputation: 31899.081657980838



Variance before imputation: 601.6565750339723

Variance after imputation: 196871.45803201268



# Arbitrary Number Imputer (arbitrary number = -1)

## Arbitrary imputation (arbitrary number = 99)

```
1 imputer_arbitrary_nn = ArbitraryNumberImputer(variables = need_impute_num, arbitrary_number = 99)
2
3 imputer_arbitrary_nn.fit(X_train, y_train)
4
5 X_train_impute_arb_nn = imputer_arbitrary_nn.transform(X_train)
6
7 X_test_impute_arb_nn = imputer_arbitrary_nn.transform(X_test)

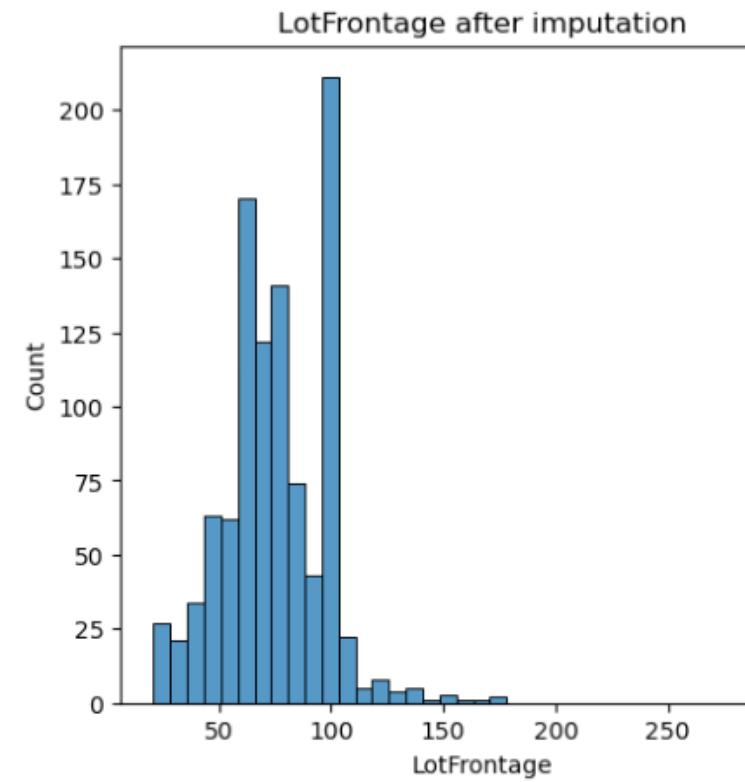
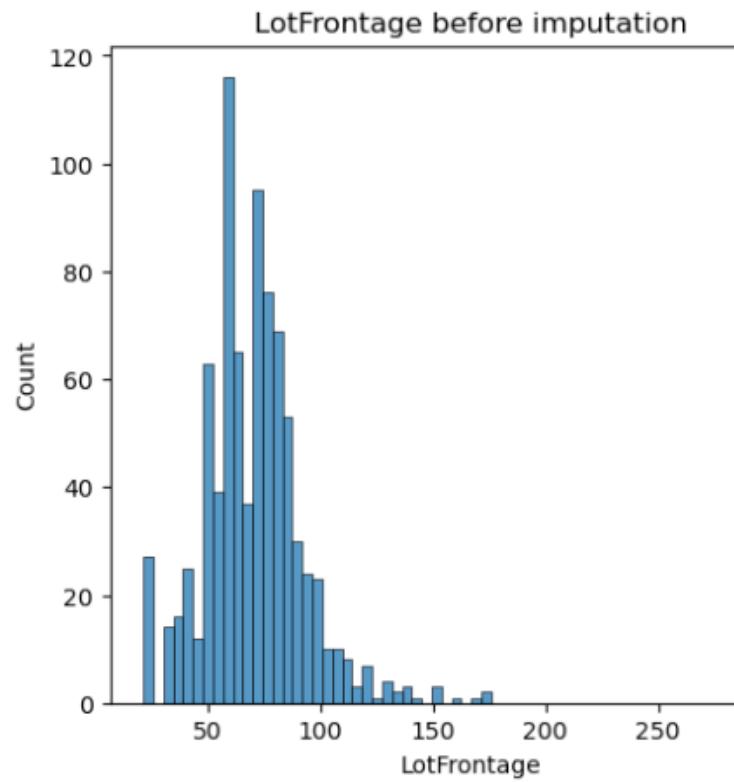
1 # Check distribution in train set after arbitrary number imputation (arbitrary value = 99)
2
3 for var in need_impute_num:
4     check_dist_and_variance(X_train, X_train_impute_arb_nn, var)
```

# Arbitrary Imputation (arbitrary number = 99)

## Check Distribution in train set after arbitrary number imputation (arbitrary value = 99)

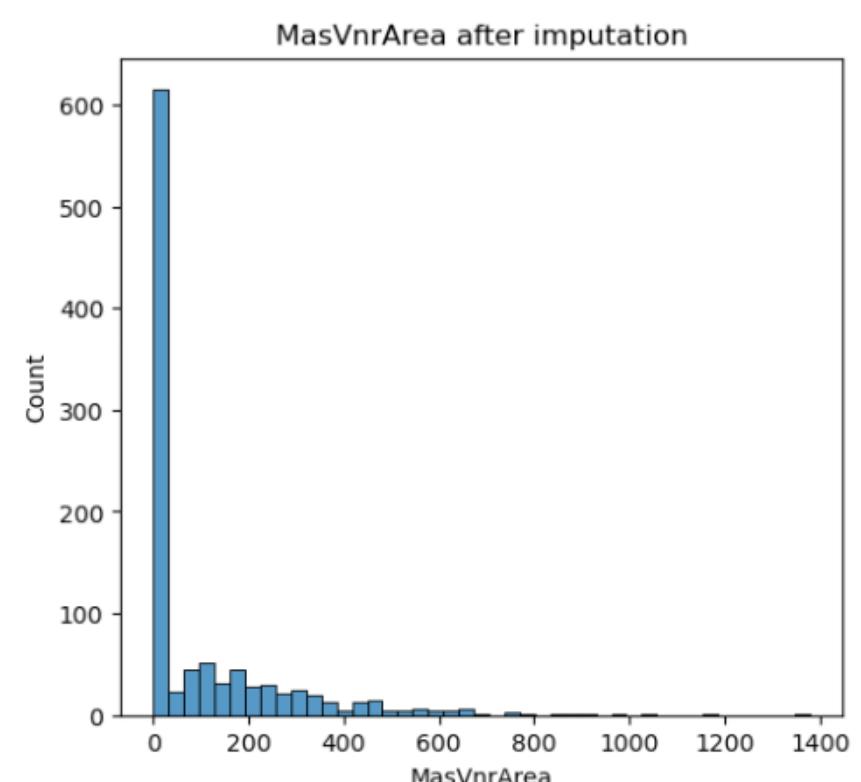
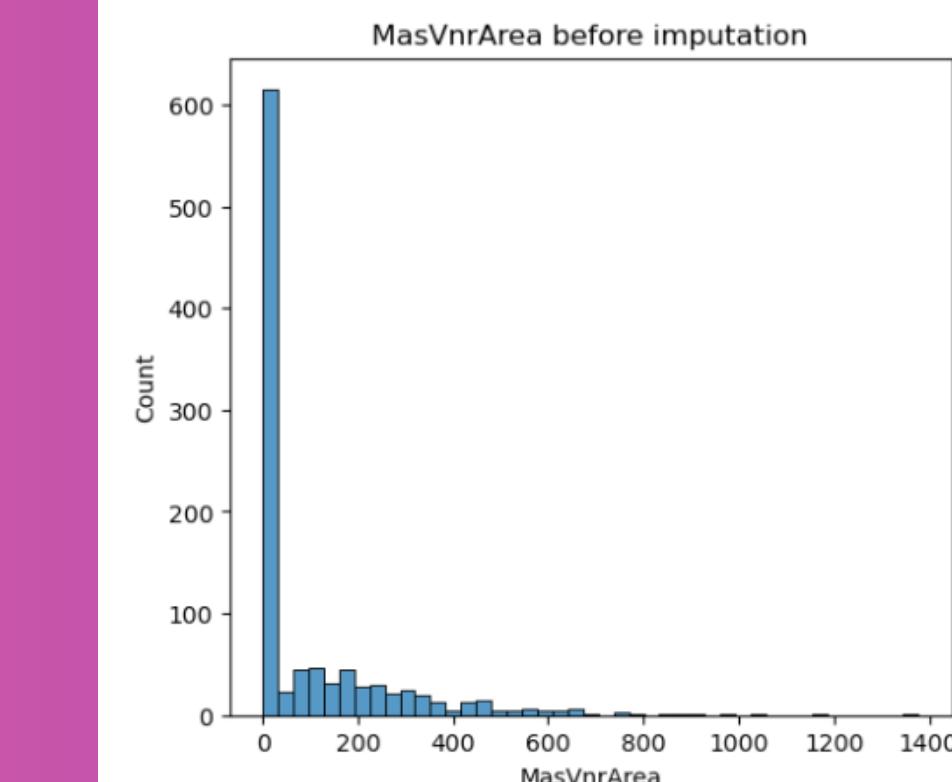
Variance before imputation: 631.2577620805455

Variance after imputation: 639.4250379985089



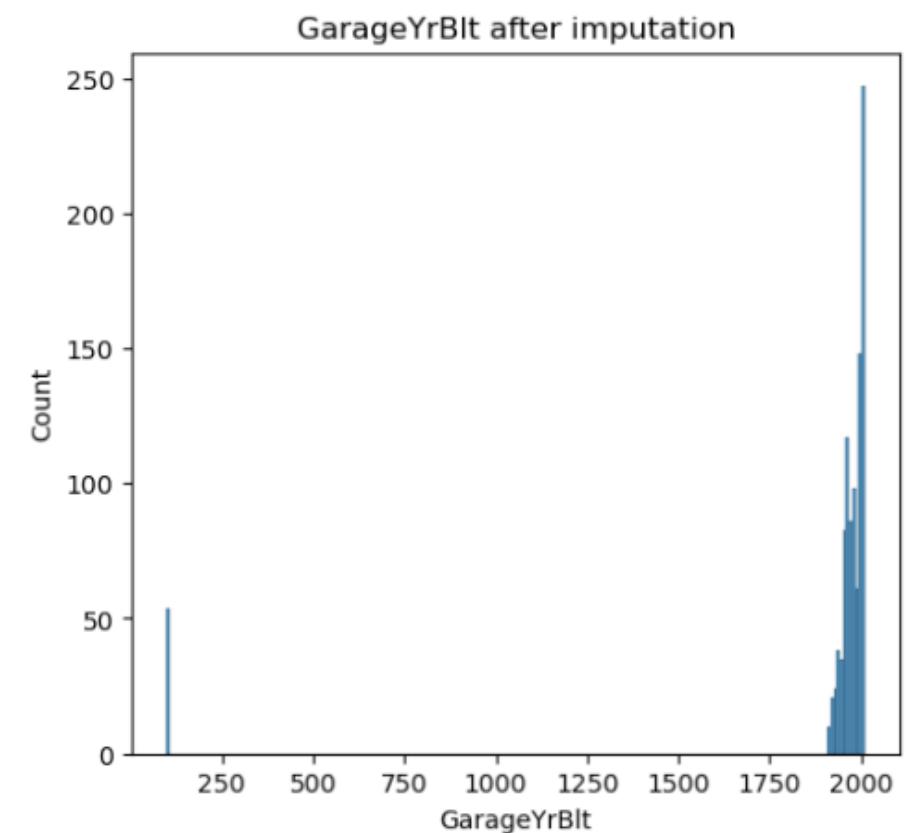
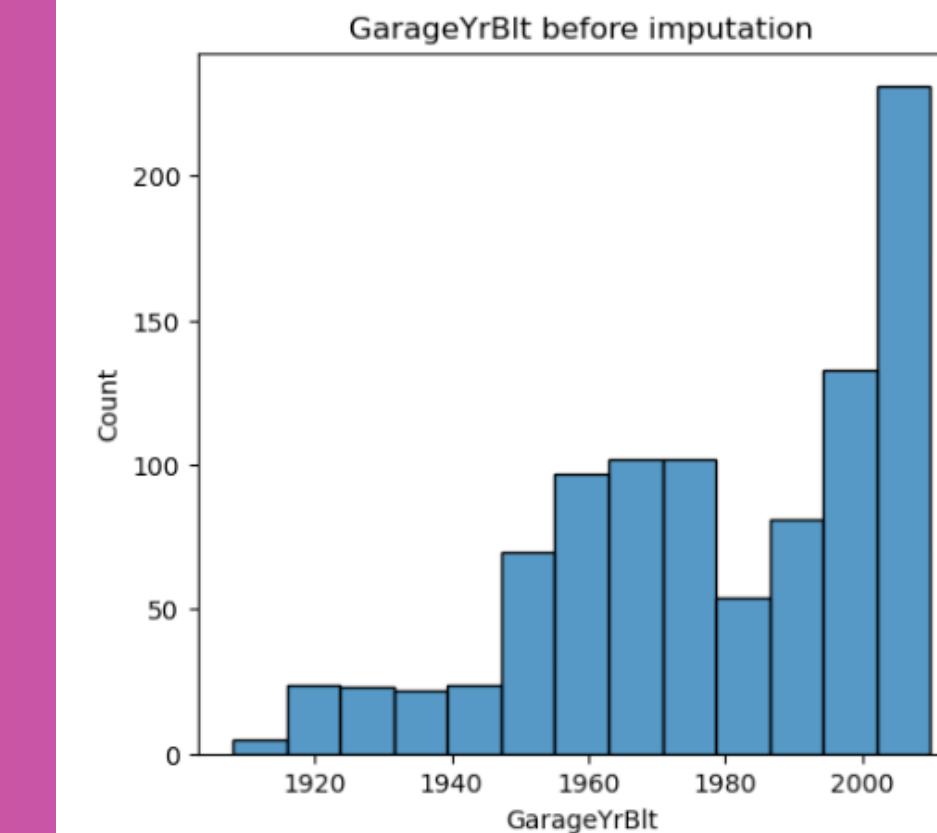
Variance before imputation: 31981.573954946656

Variance after imputation: 31856.366158997647



# Arbitrary Imputation (arbitrary number = 99)

Variance after imputation: 177539.4164243643



## Mean&Median imputation

### Median imputation

```
1 imputer_median = MeanMedianImputer(variables = need_impute_num, imputation_method = 'median')
2
3 imputer_median.fit(X_train, y_train)
4
5 X_train_impute_median = imputer_median.transform(X_train)
6
7 X_test_impute_median = imputer_median.transform(X_test)
```

```
1 # Distribution after median imputation
2
3 for var in need_impute_num:
4
5     check_dist_and_variance(X_train, X_train_impute_median, var)
```

# Mean Median Imputation:

## Median Imputation

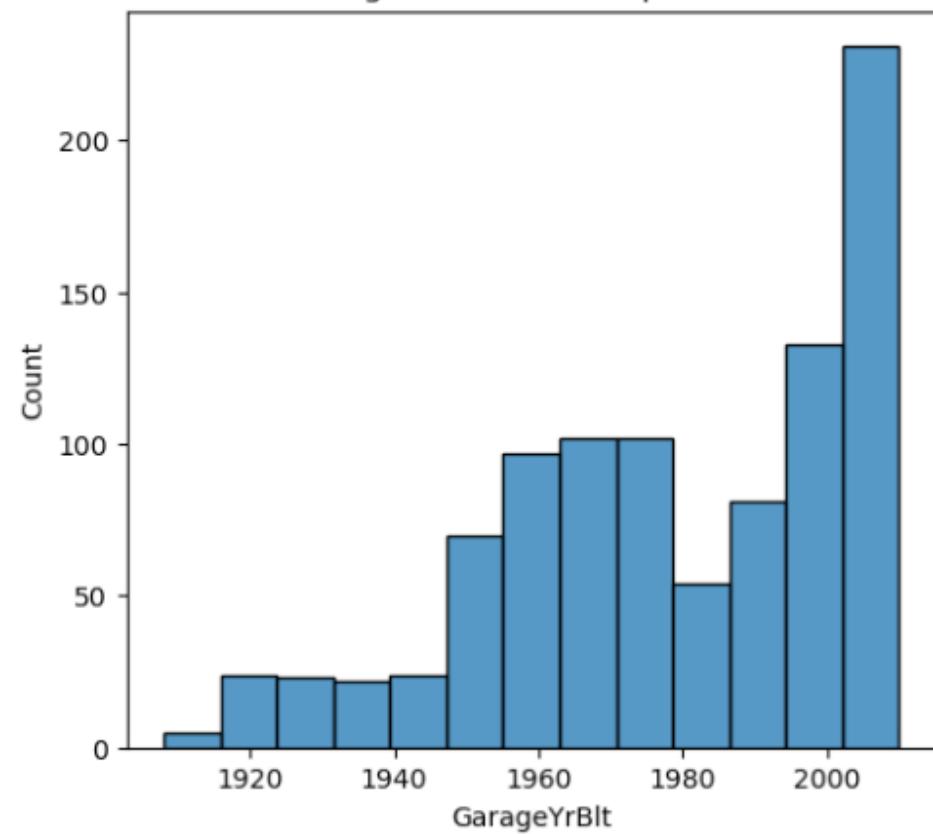
### Distribution after median imputation

# Distribution after median imputation

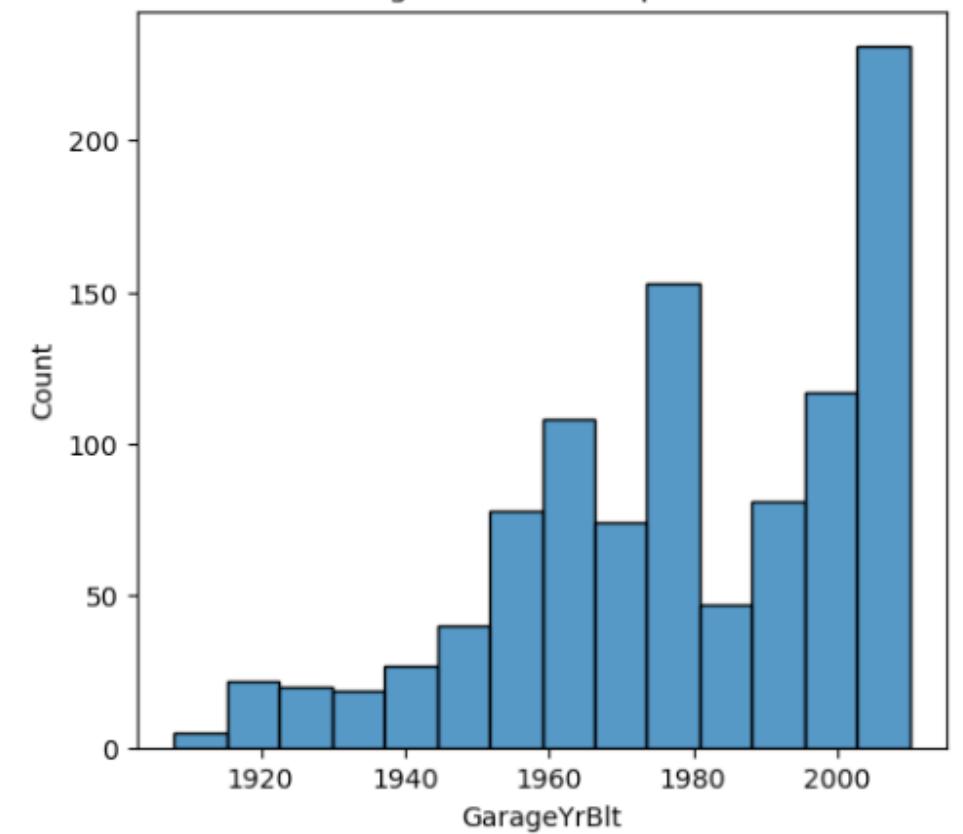
Variance before imputation: 601.6565750339723

Variance after imputation: 569.9412005420418

GarageYrBlt before imputation



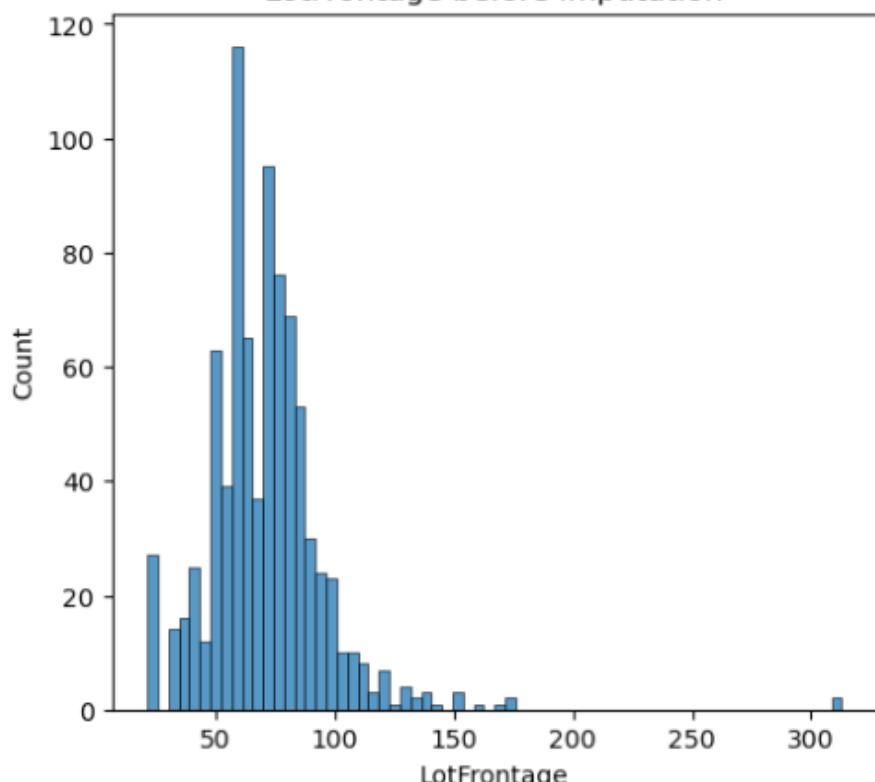
GarageYrBlt after imputation



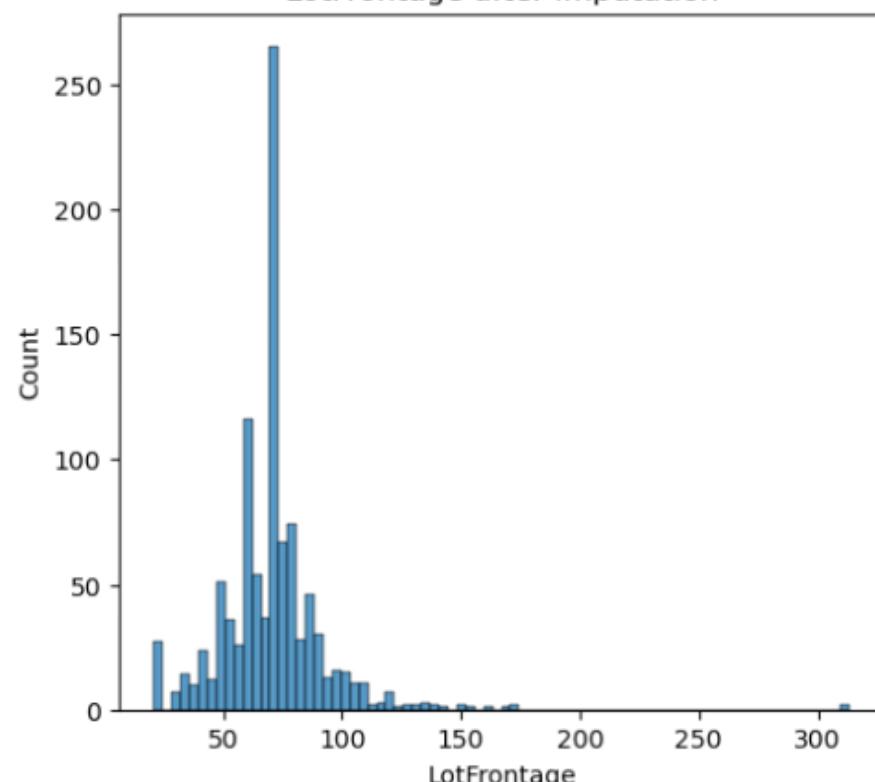
Variance before imputation: 631.2577620805455

Variance after imputation: 519.9834867010011

LotFrontage before imputation



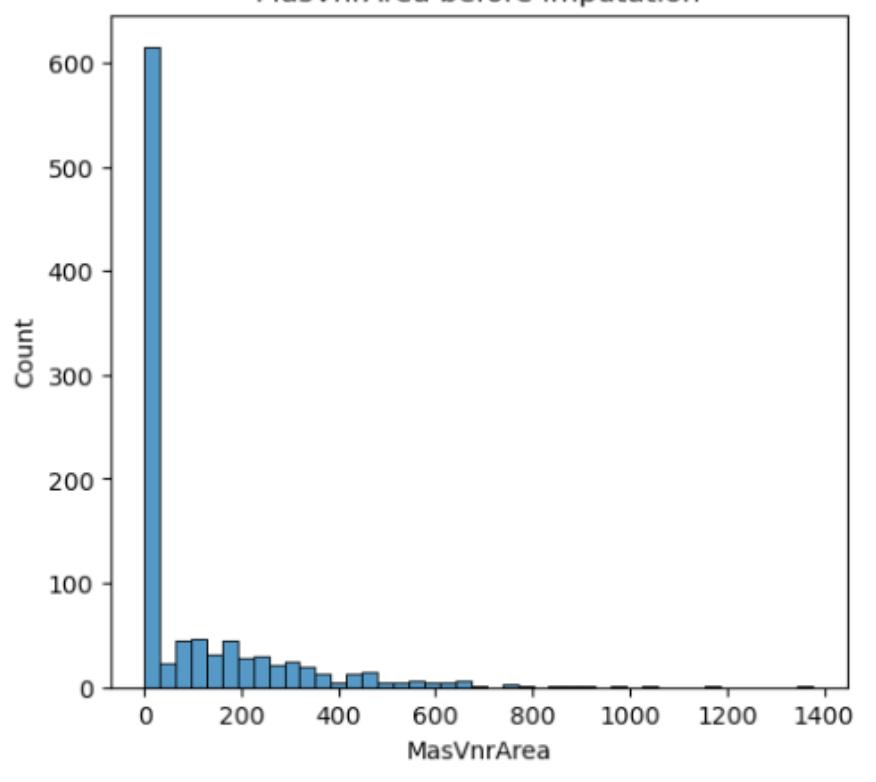
LotFrontage after imputation



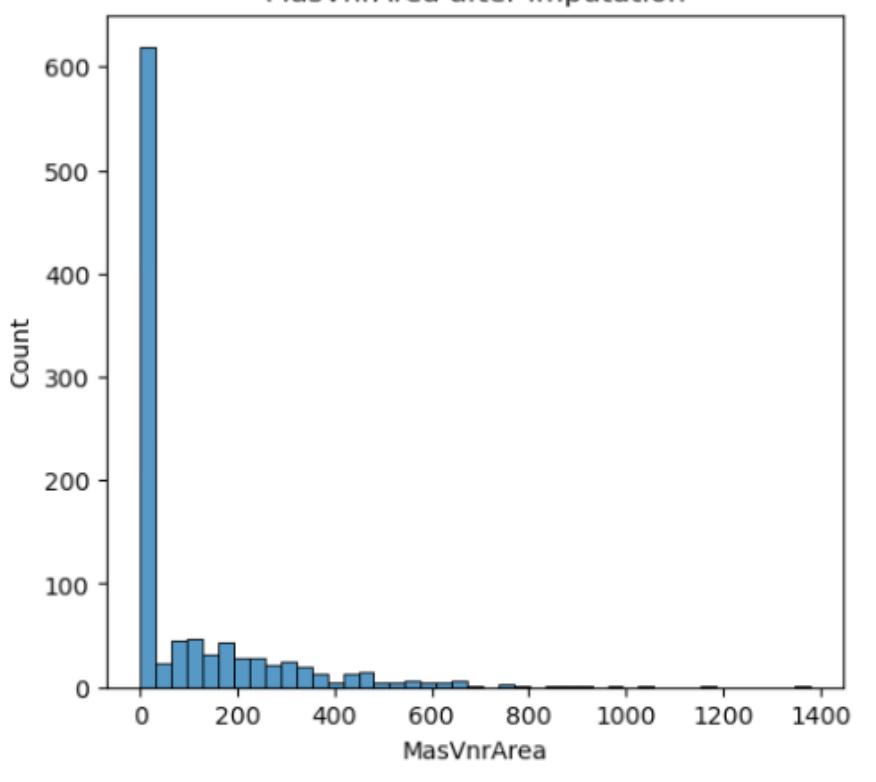
Variance before imputation: 31981.573954946656

Variance after imputation: 31898.268165970585

MasVnrArea before imputation



MasVnrArea after imputation



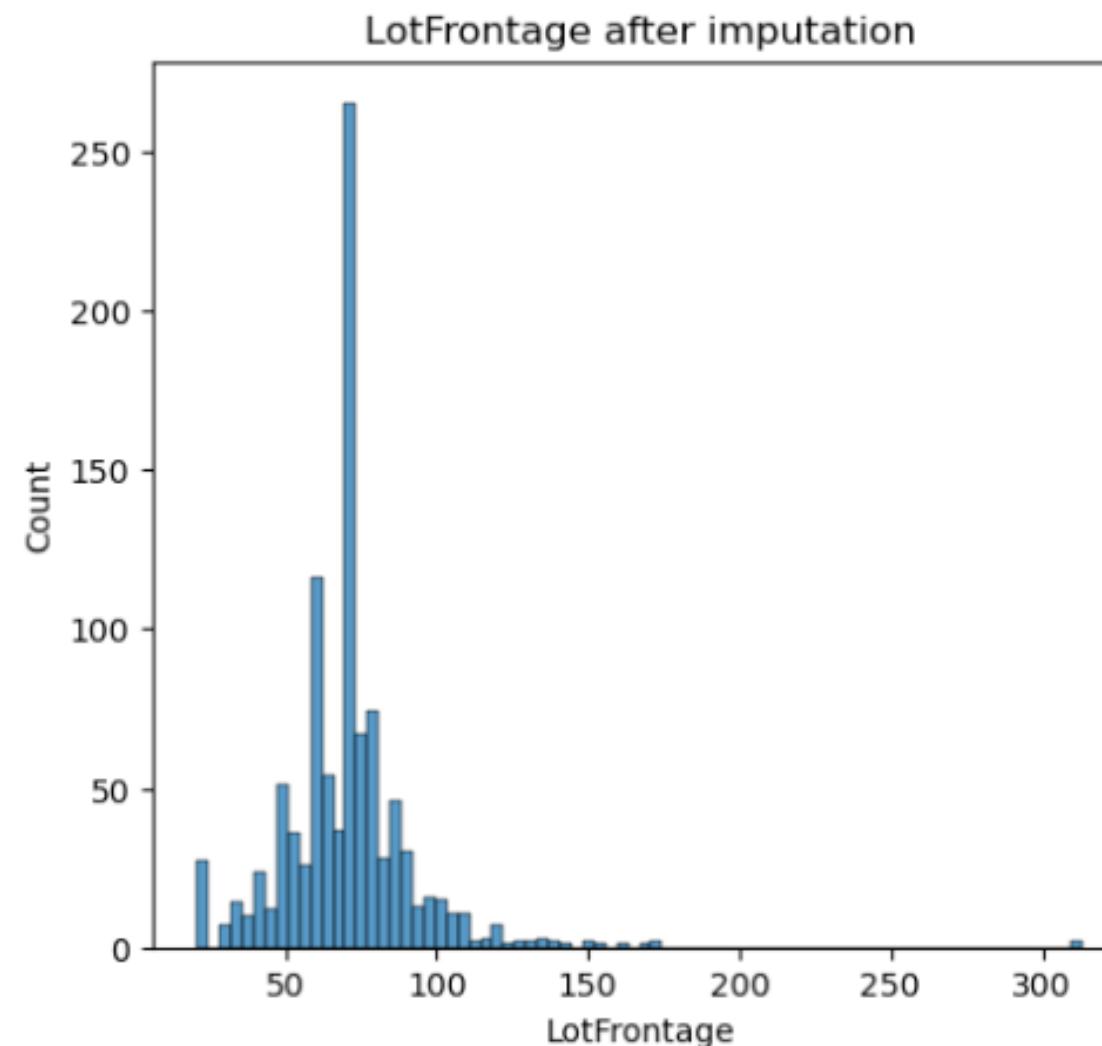
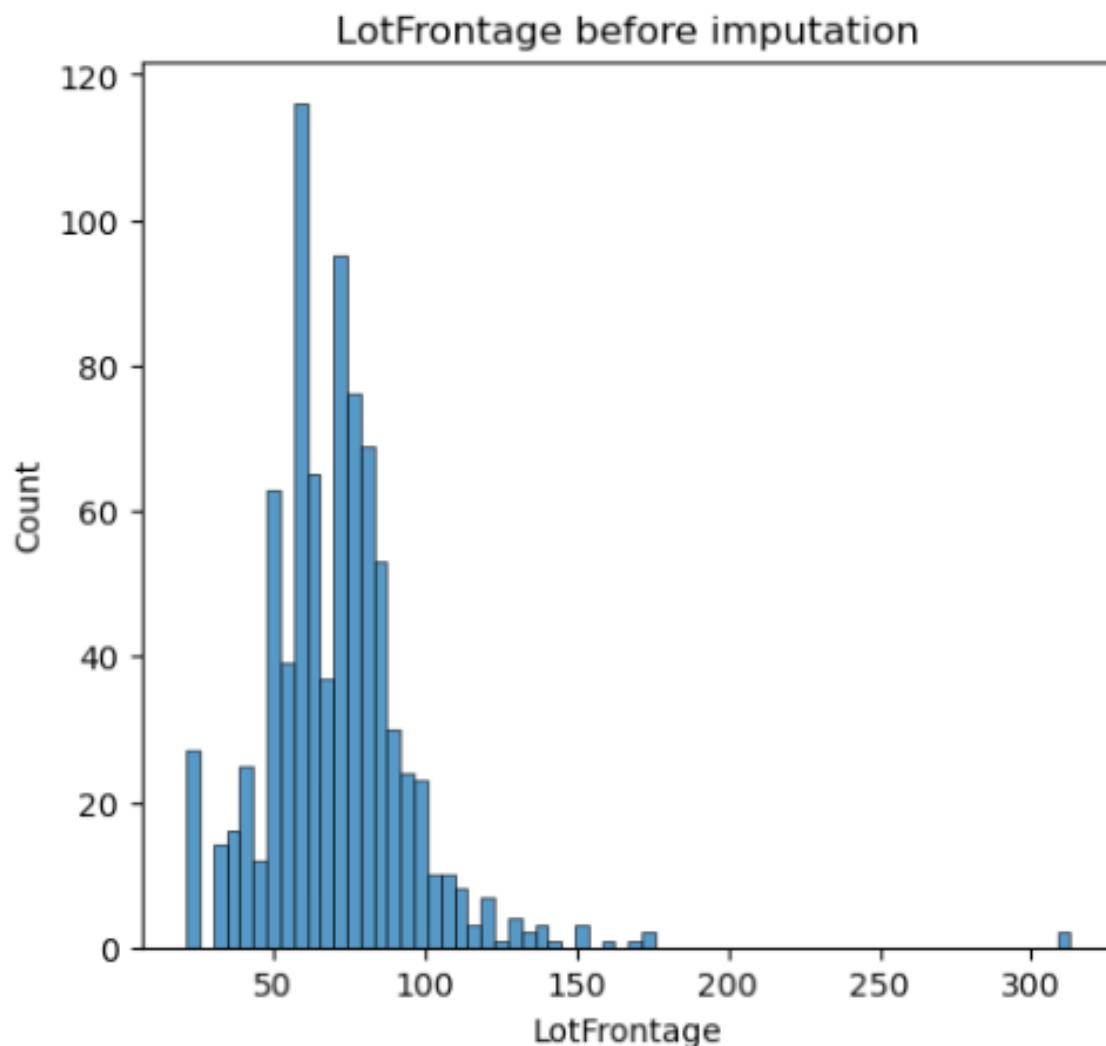
## Mean imputation

```
1 imputer_mean = MeanMedianImputer(variables = ['LotFrontage'], imputation_method = 'median')
2
3 imputer_mean.fit(X_train, y_train)
4
5 X_train_impute_mean = imputer_mean.transform(X_train)
6
7 X_test_impute_mean = imputer_mean.transform(X_test)
```

```
1 check_dist_and_variance(X_train, X_train_impute_mean, 'LotFrontage')
```

Variance before imputation: 631.2577620805455

Variance after imputation: 519.9834867010011



**Check Distribution after  
mean imputation**

**mean imputation  
method**

## Random sample imputation

```
1 imputer_rs = RandomSampleImputer(variables = need_impute, random_state = 4)
2
3 imputer_rs.fit(X_train)
4
5 X_train_impute_rs = imputer_rs.transform(X_train)
6
7 X_test_impute_rs = imputer_rs.transform(X_test)

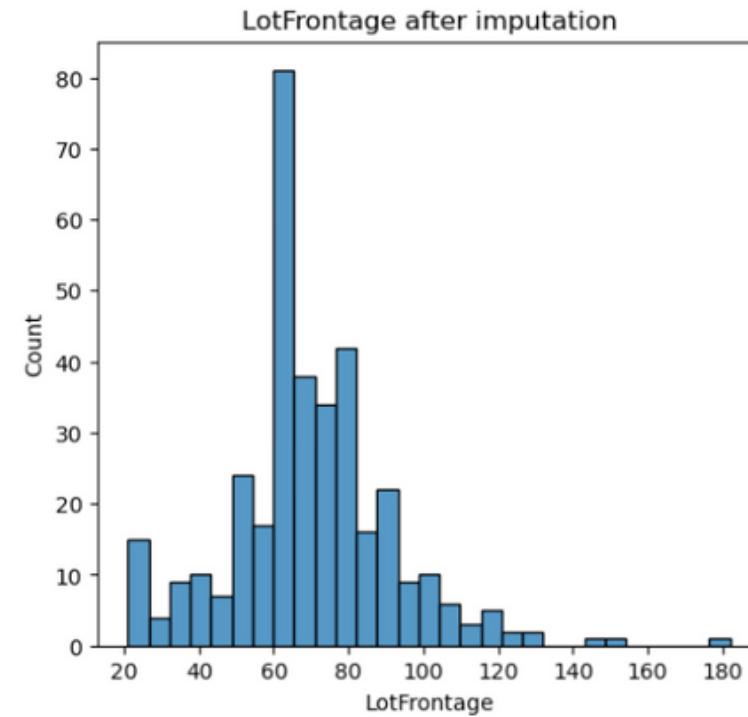
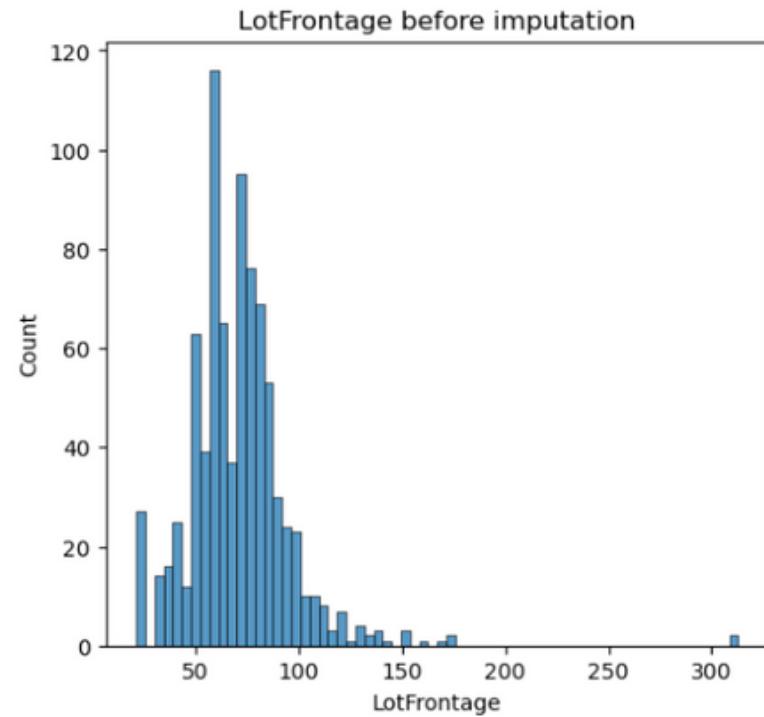
1 # Distribution after random sample imputer
2 for var in need_impute:
3     check_dist_and_variance(X_train, X_test, var)
```

# Random Sample Imputation

## Distribution after Random Sample Imputation

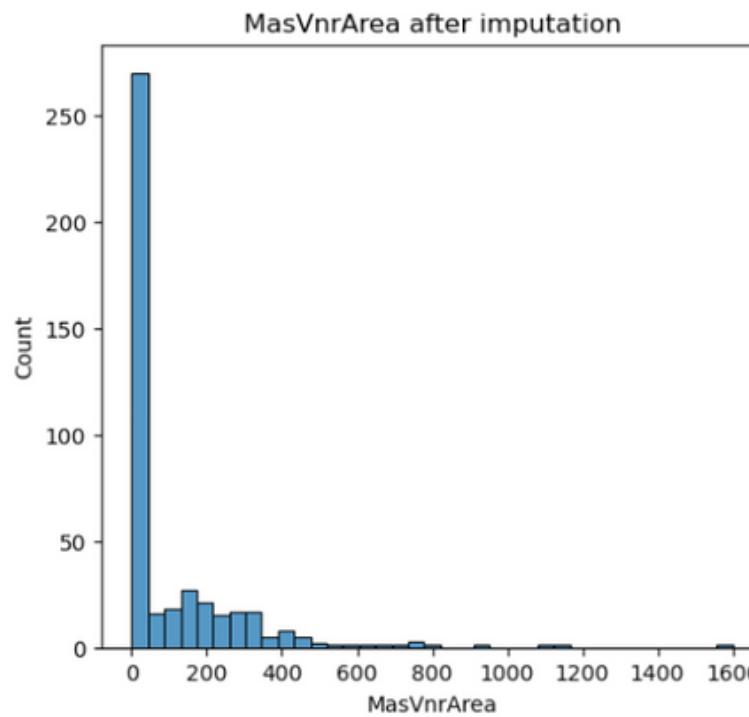
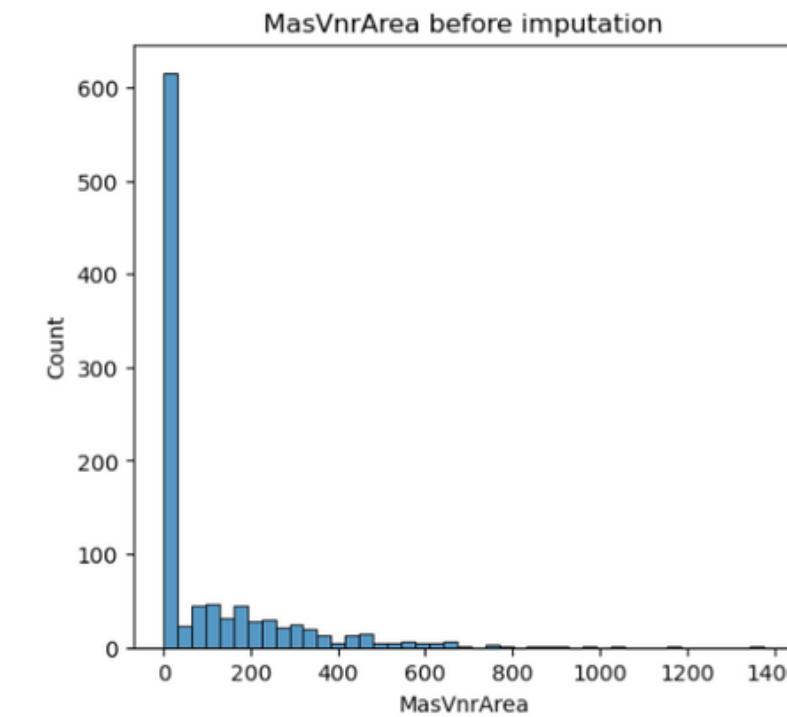
Variance before imputation: 631.2577620805455

Variance after imputation: 493.3044303698978

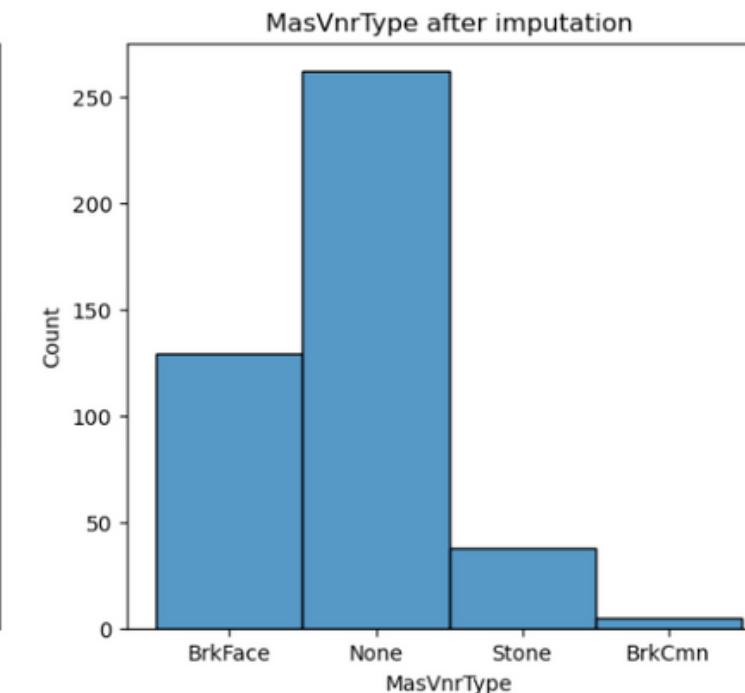
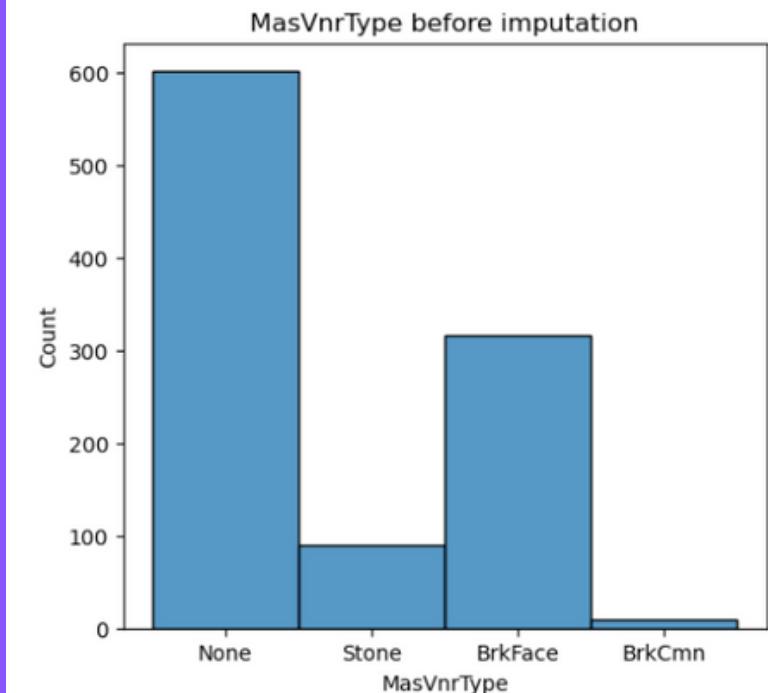


Variance before imputation: 31981.573954946656

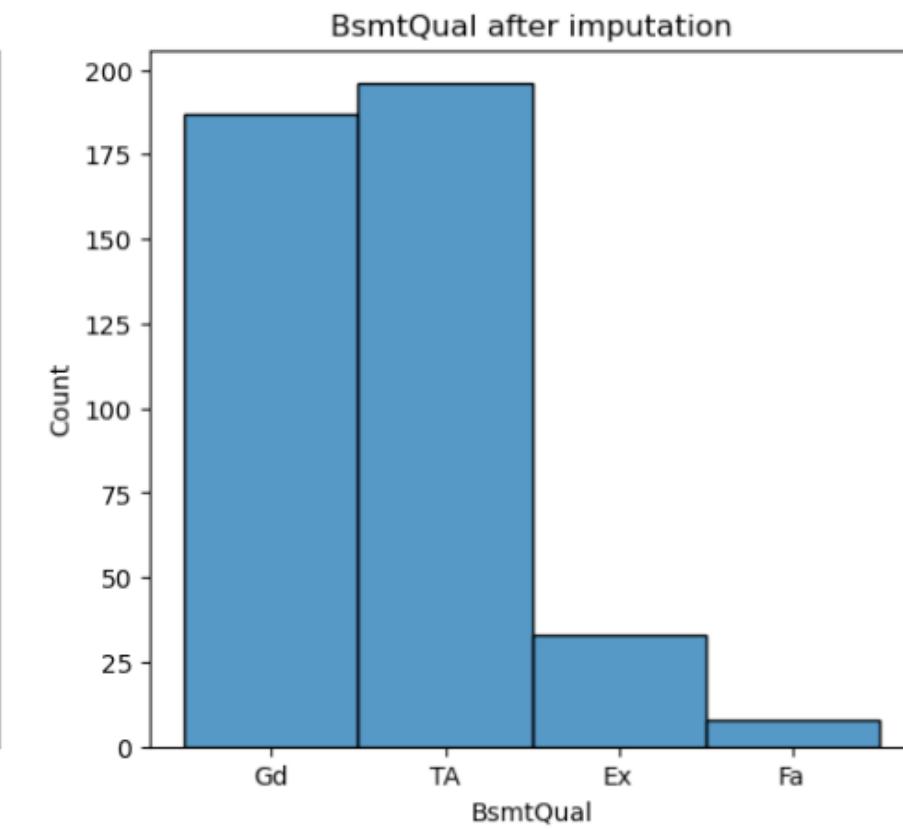
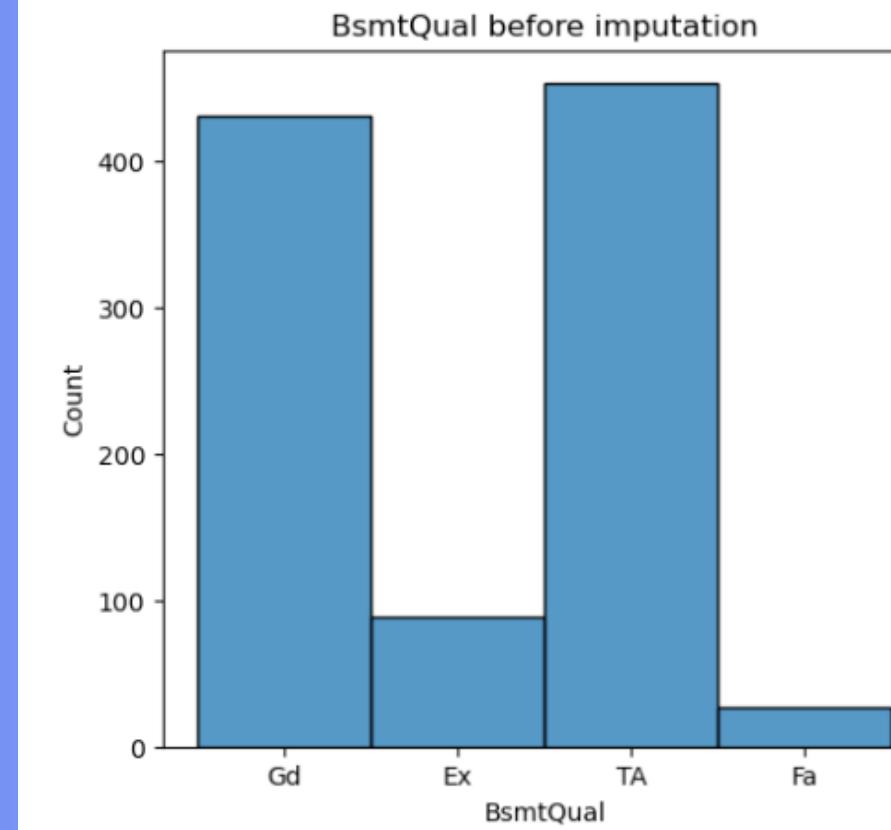
Variance after imputation: 34747.63435893616



Categorical variable



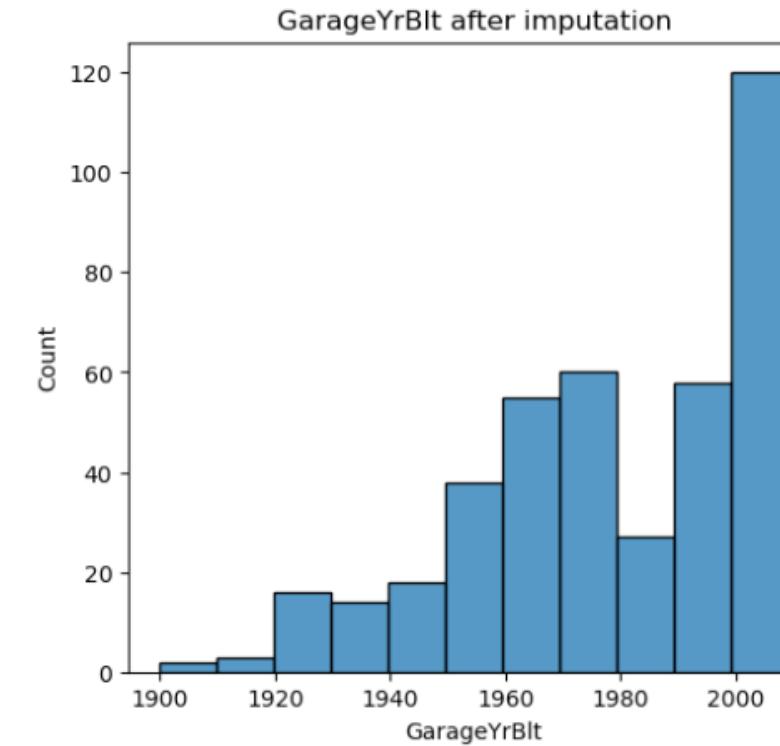
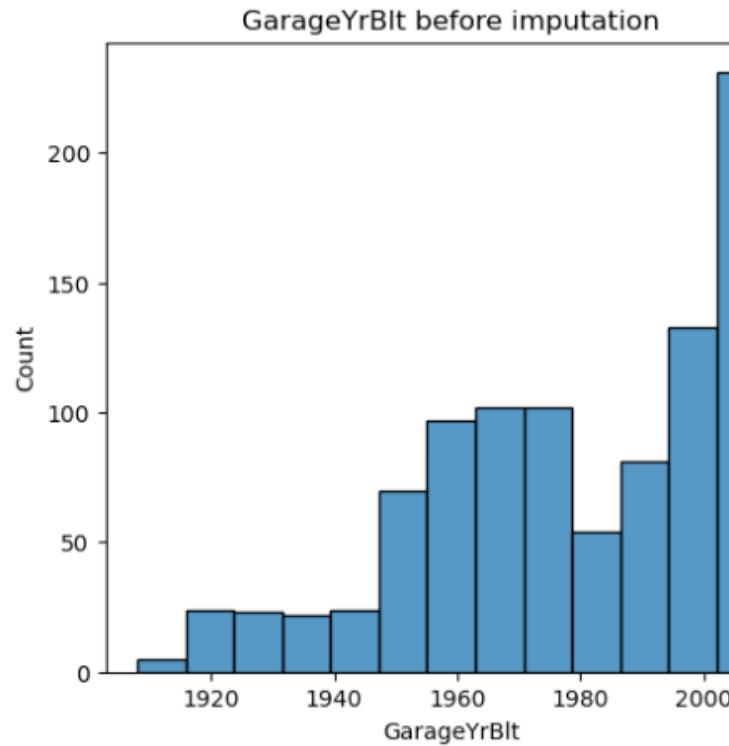
Categorical variable



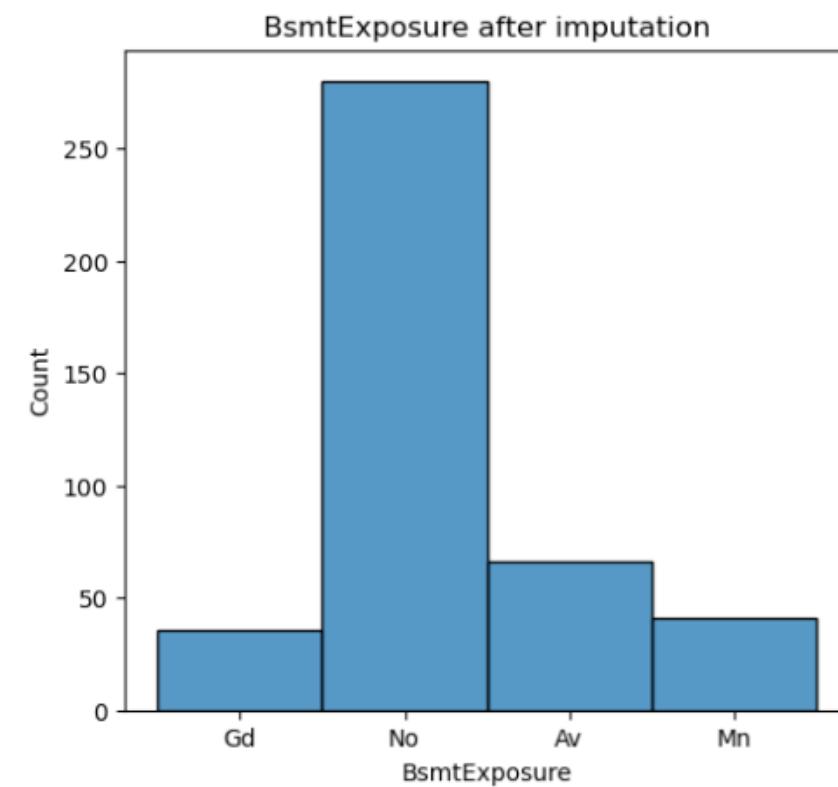
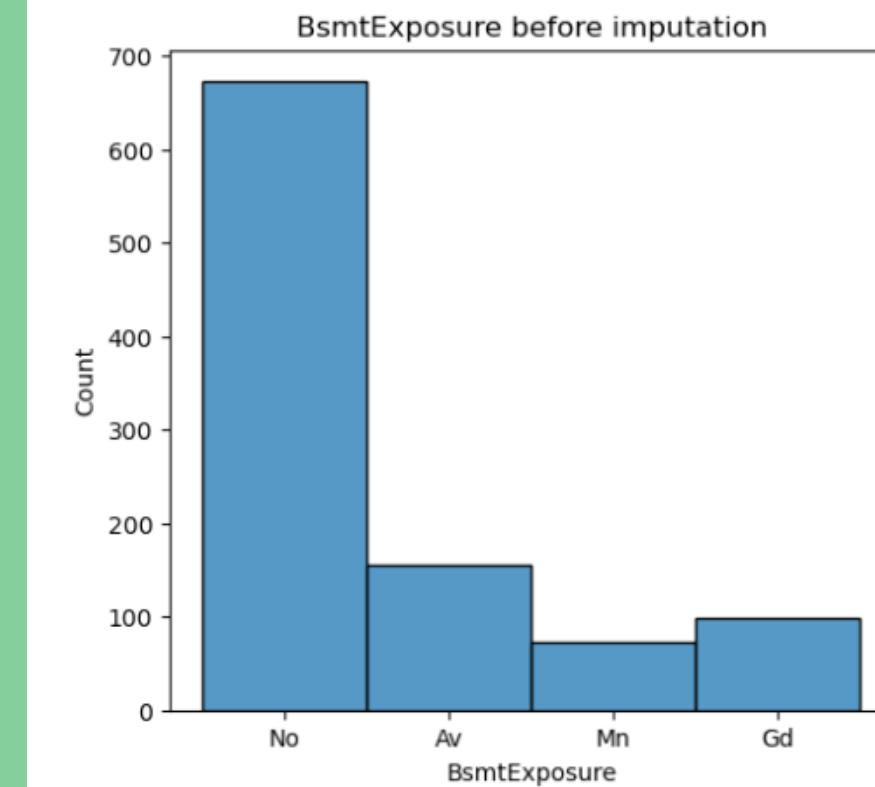
# Distribution after Random Sample Imputation

Variance before imputation: 601.6565750339723

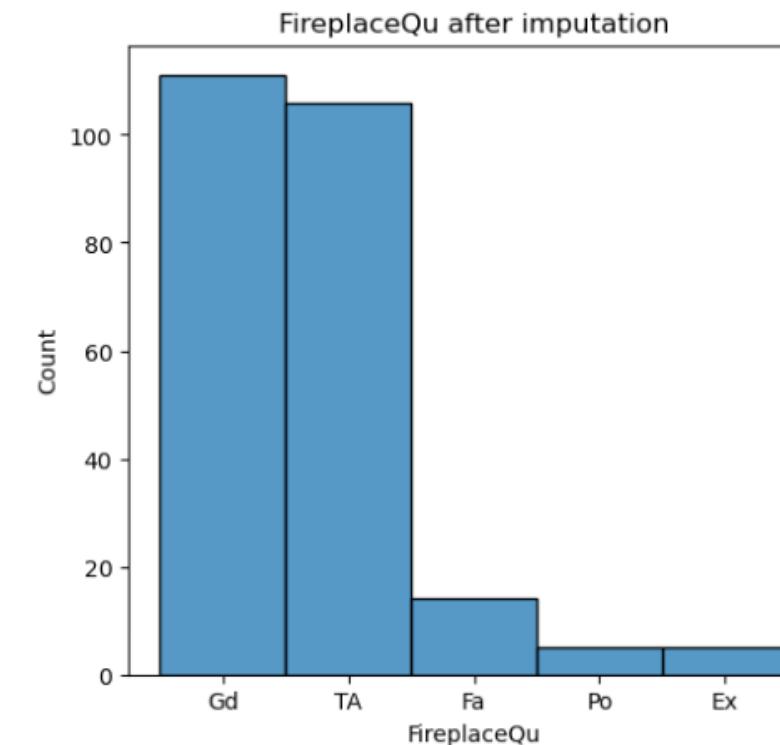
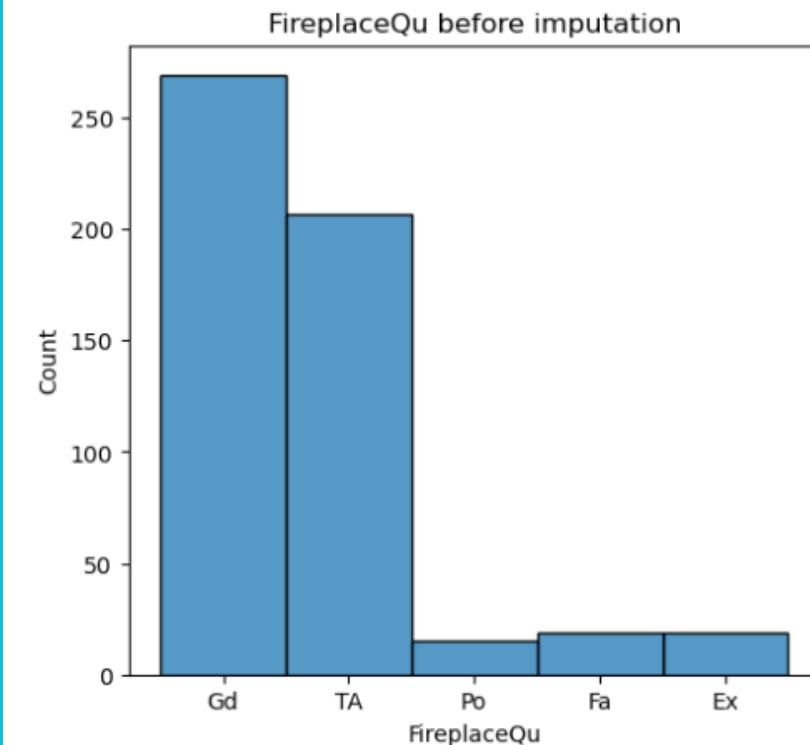
Variance after imputation: 629.7500207702809



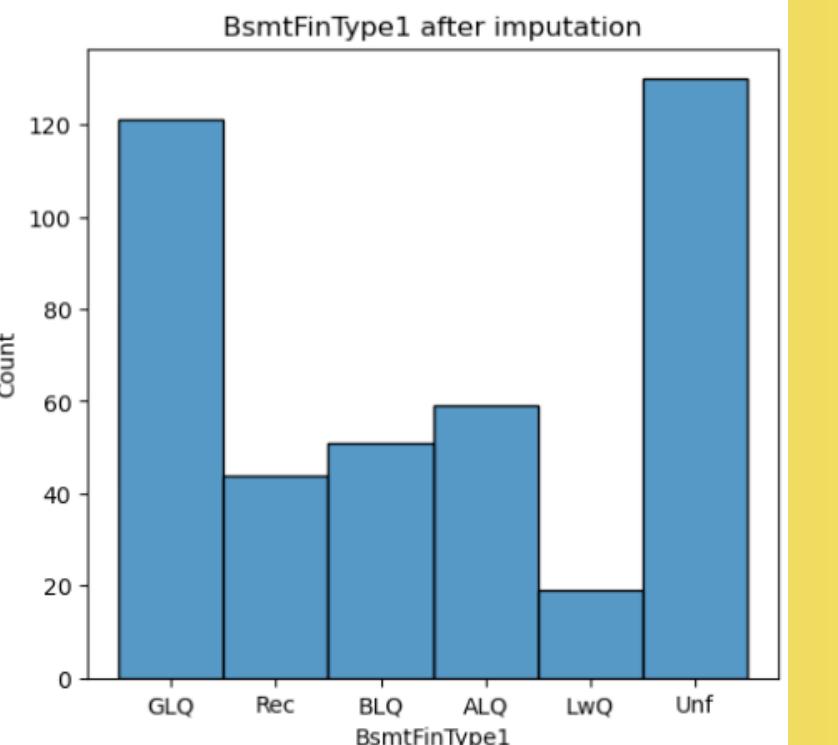
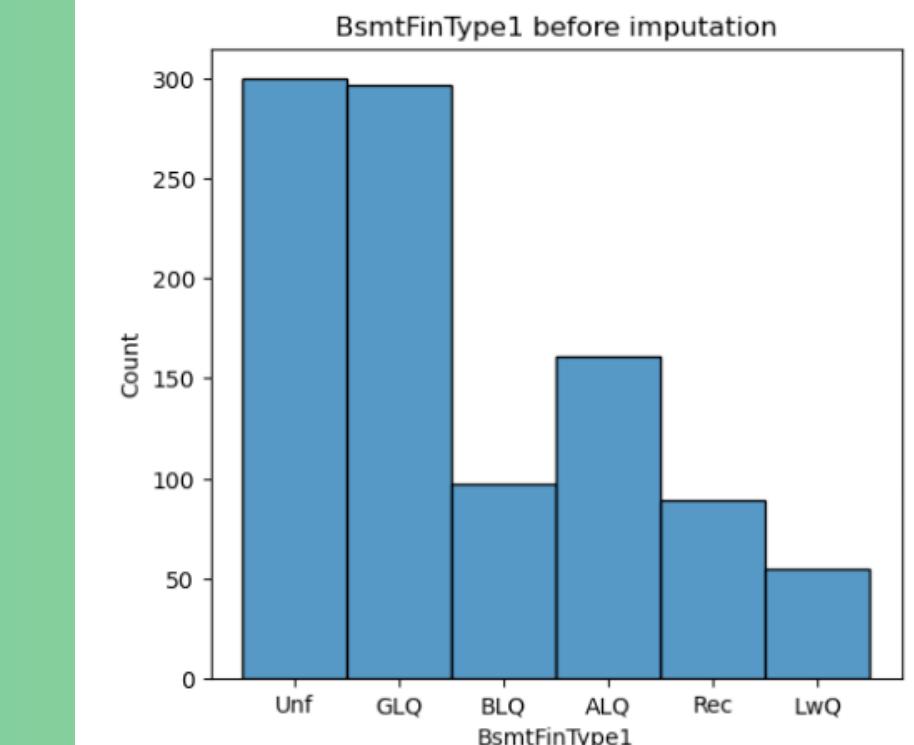
Categorical variable



Categorical variable



Categorical variable



# Distribution after Random Sample Imputation

## Final imputation stack

```
1 pipe_imputation = Pipeline(steps = [
2     ('frequent', CategoricalImputer(
3         variables = ['MasVnrType', 'BsmtCond', 'BsmtExposure', 'BsmtFinType2', 'Electrical', 'BsmtFinType1'],
4         imputation_method = 'frequent')
5     ),
6     ('missng_labels', CategoricalImputer(
7         variables = ['FireplaceQu', 'BsmtQual', 'GarageType', 'GarageQual', 'GarageCond', 'GarageFinish'],
8         fill_value = 'missing',
9         imputation_method = 'missing')
10    ),
11    ('median', MeanMedianImputer(
12        variables = ['MasVnrArea'],
13        imputation_method = 'median'
14    )
15    ),
16    ('random_sample', RandomSampleImputer(
17        variables = ['LotFrontage', 'GarageYrBlt'],
18        random_state = 4
19    )
20    )
21 ])
22
23 pipe_imputation.fit(X_train)
```

# Final Imputation Stack with PipeLine

```
1 X_train_impute = pipe_imputation.transform(X_train)
2
3 X_test_impute = pipe_imputation.transform(X_test)
```

```
1 # Check imputation features in sets
2 print('Train set \n')
3 print(X_train_impute[need_impute].isnull().mean())
4
5 print('\n Test set \n')
6 print(X_test_impute[need_impute].isnull().mean())
```

Train set

```
LotFrontage      0.0
MasVnrType       0.0
MasVnrArea       0.0
BsmtQual         0.0
BsmtCond         0.0
BsmtExposure     0.0
BsmtFinType1     0.0
BsmtFinType2     0.0
Electrical        0.0
FireplaceQu      0.0
GarageType        0.0
GarageYrBlt      0.0
GarageFinish      0.0
GarageQual        0.0
GarageCond        0.0
dtype: float64
```

Test set

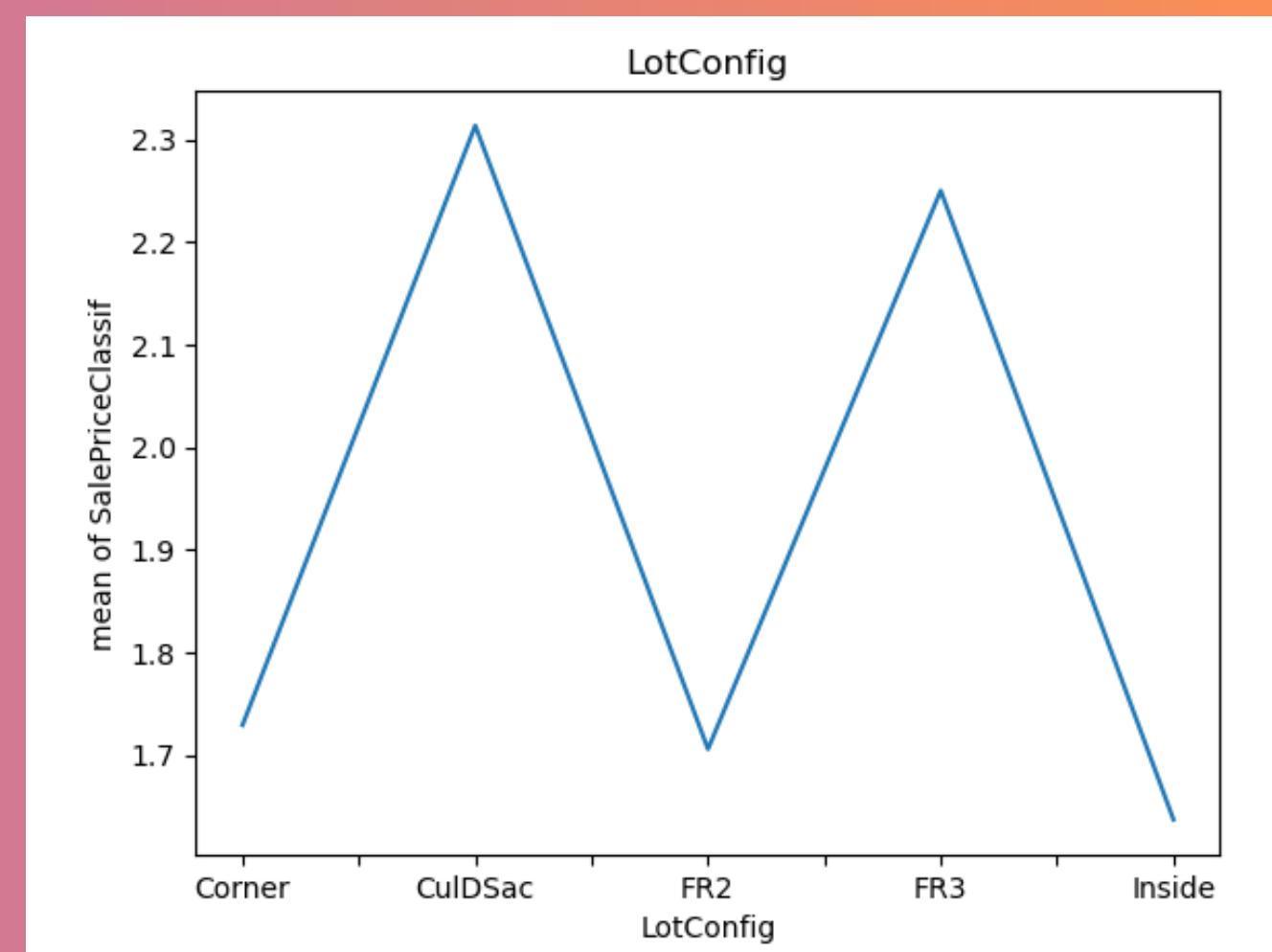
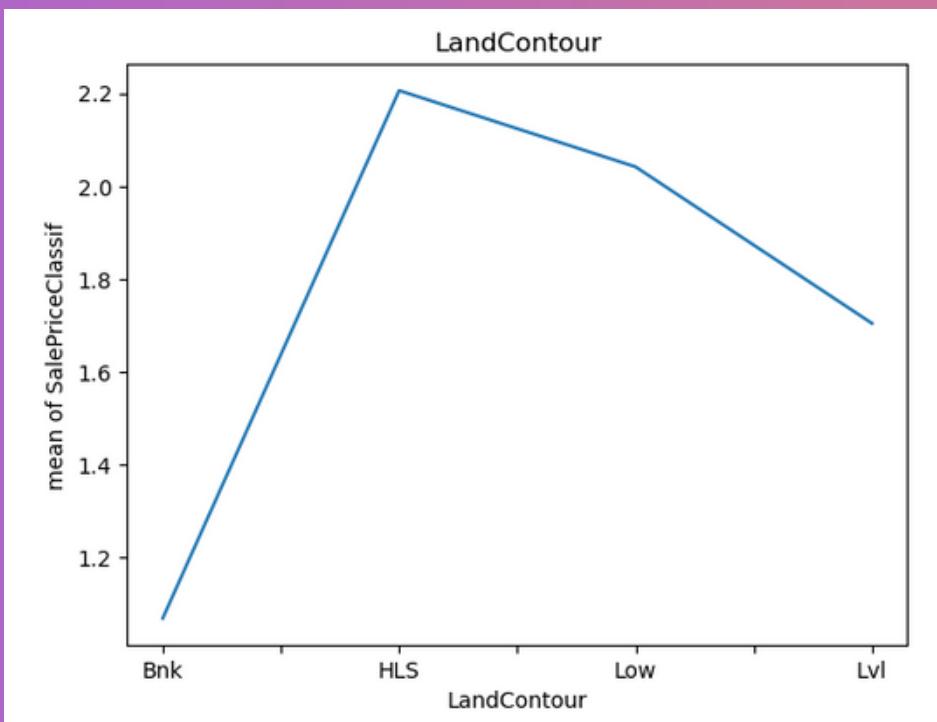
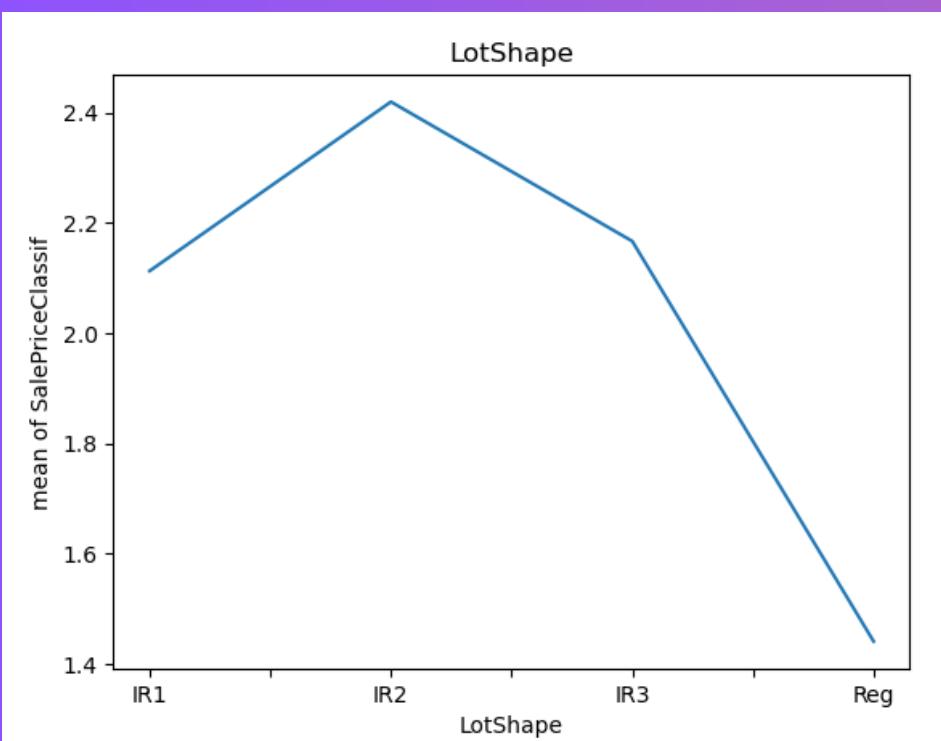
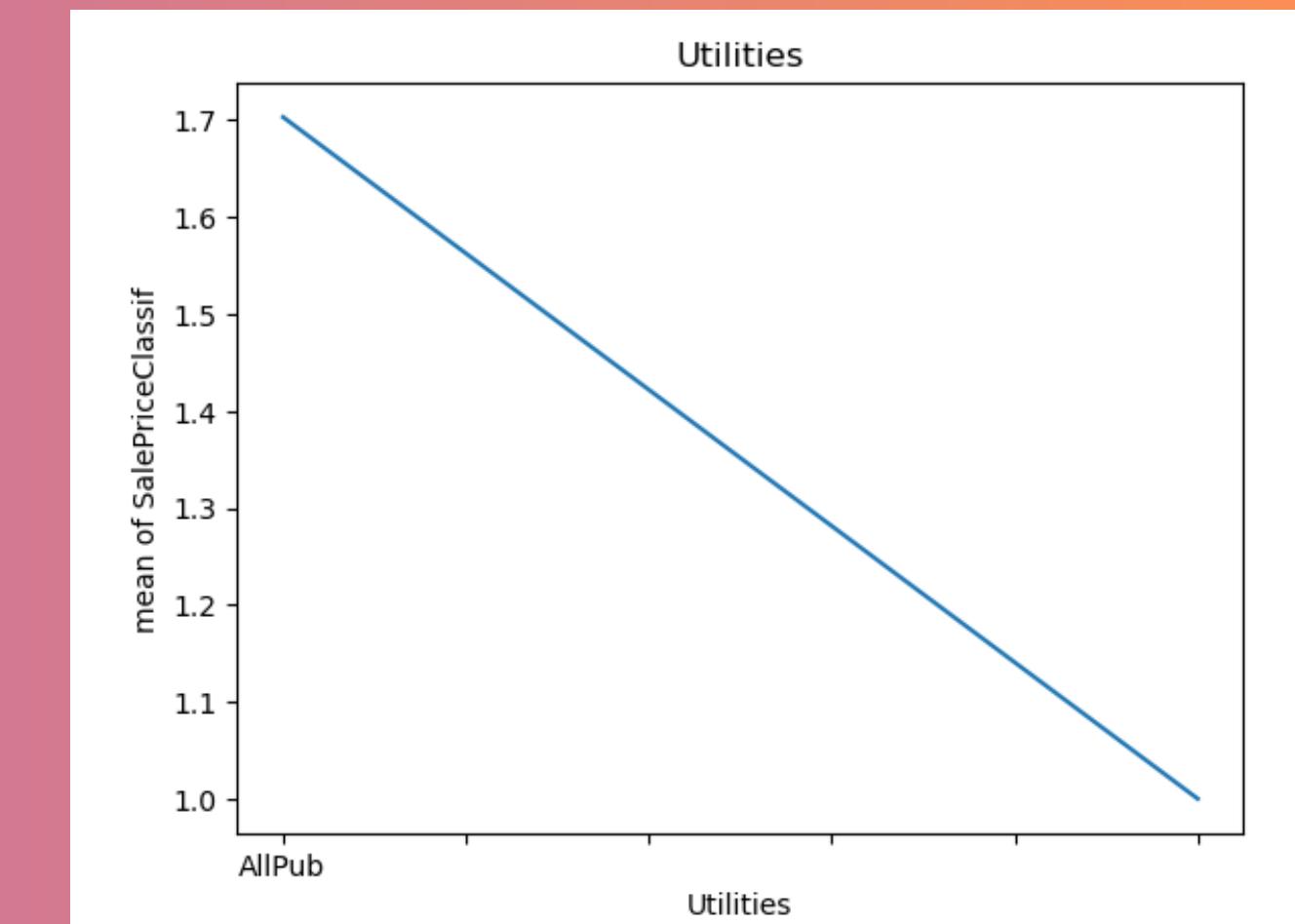
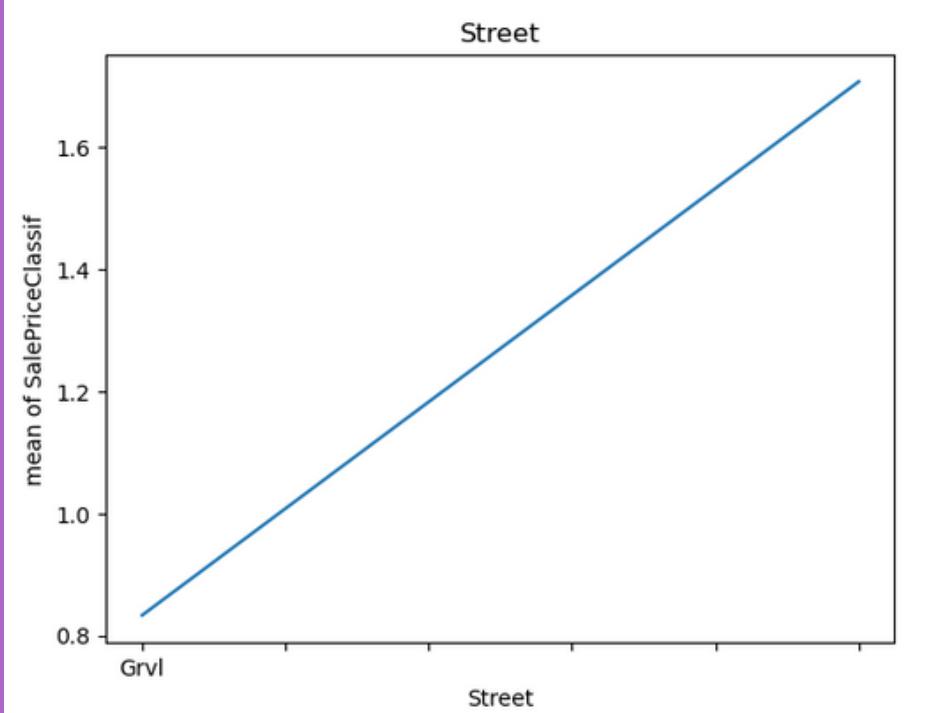
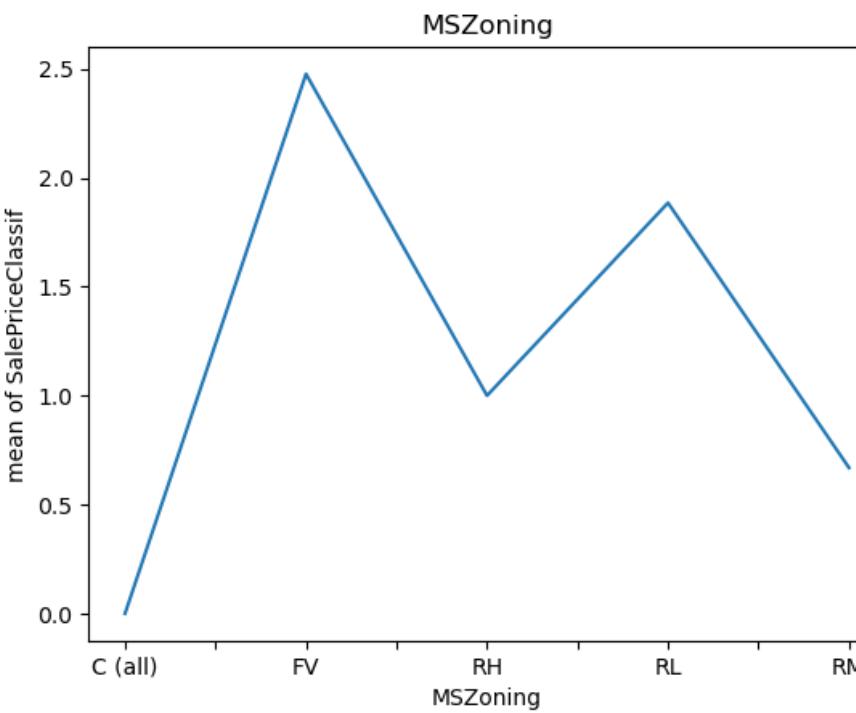
```
LotFrontage      0.0
MasVnrType       0.0
MasVnrArea       0.0
BsmtQual         0.0
BsmtCond         0.0
BsmtExposure     0.0
BsmtFinType1     0.0
BsmtFinType2     0.0
Electrical        0.0
FireplaceQu      0.0
GarageType        0.0
GarageYrBlt      0.0
GarageFinish      0.0
GarageQual        0.0
GarageCond        0.0
dtype: float64
```

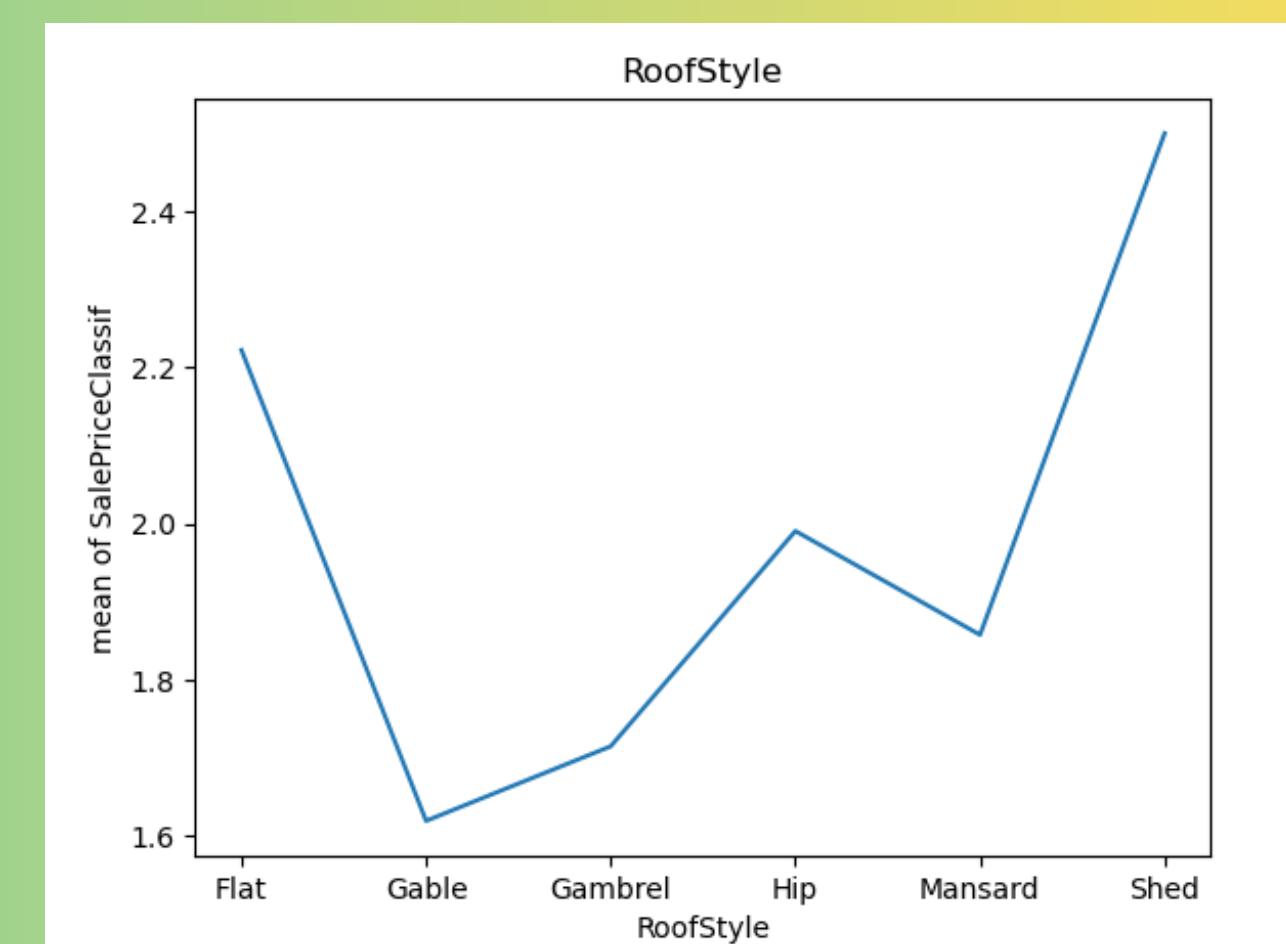
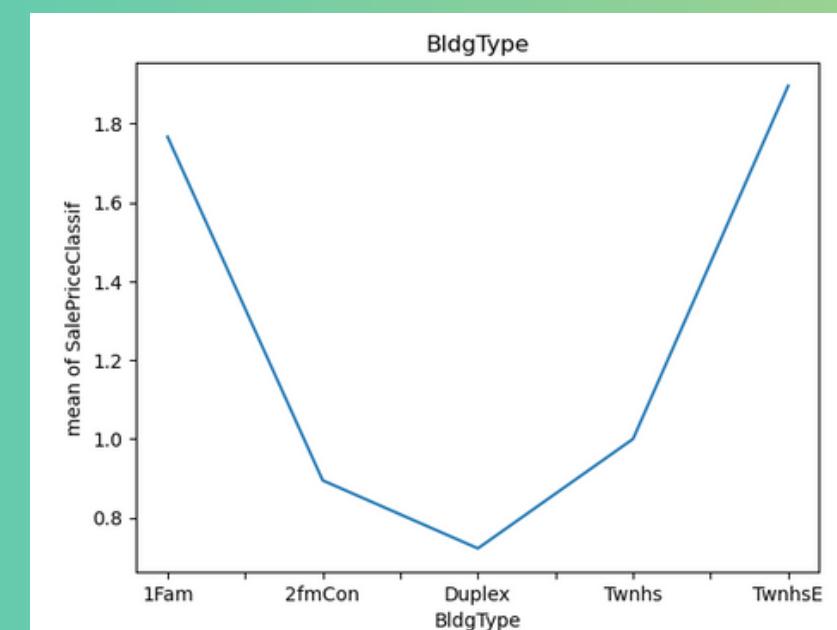
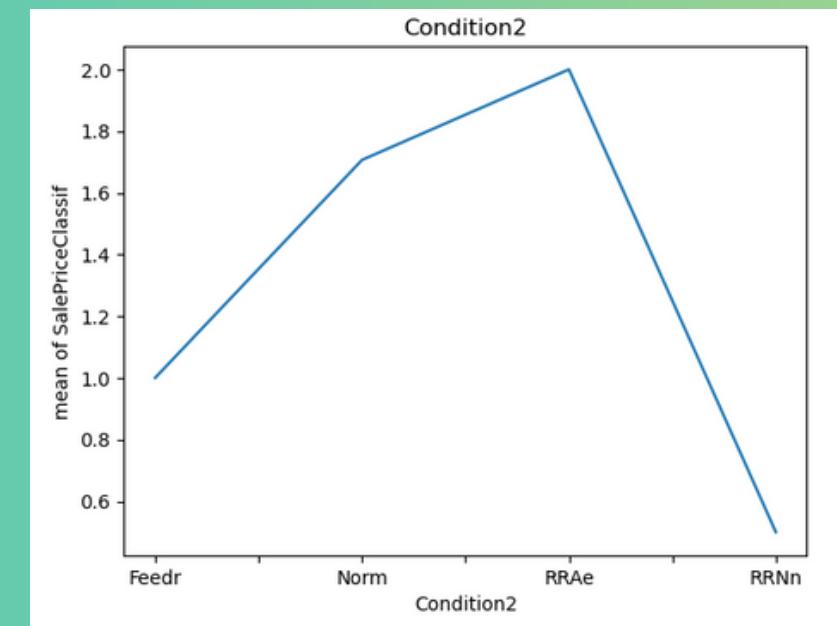
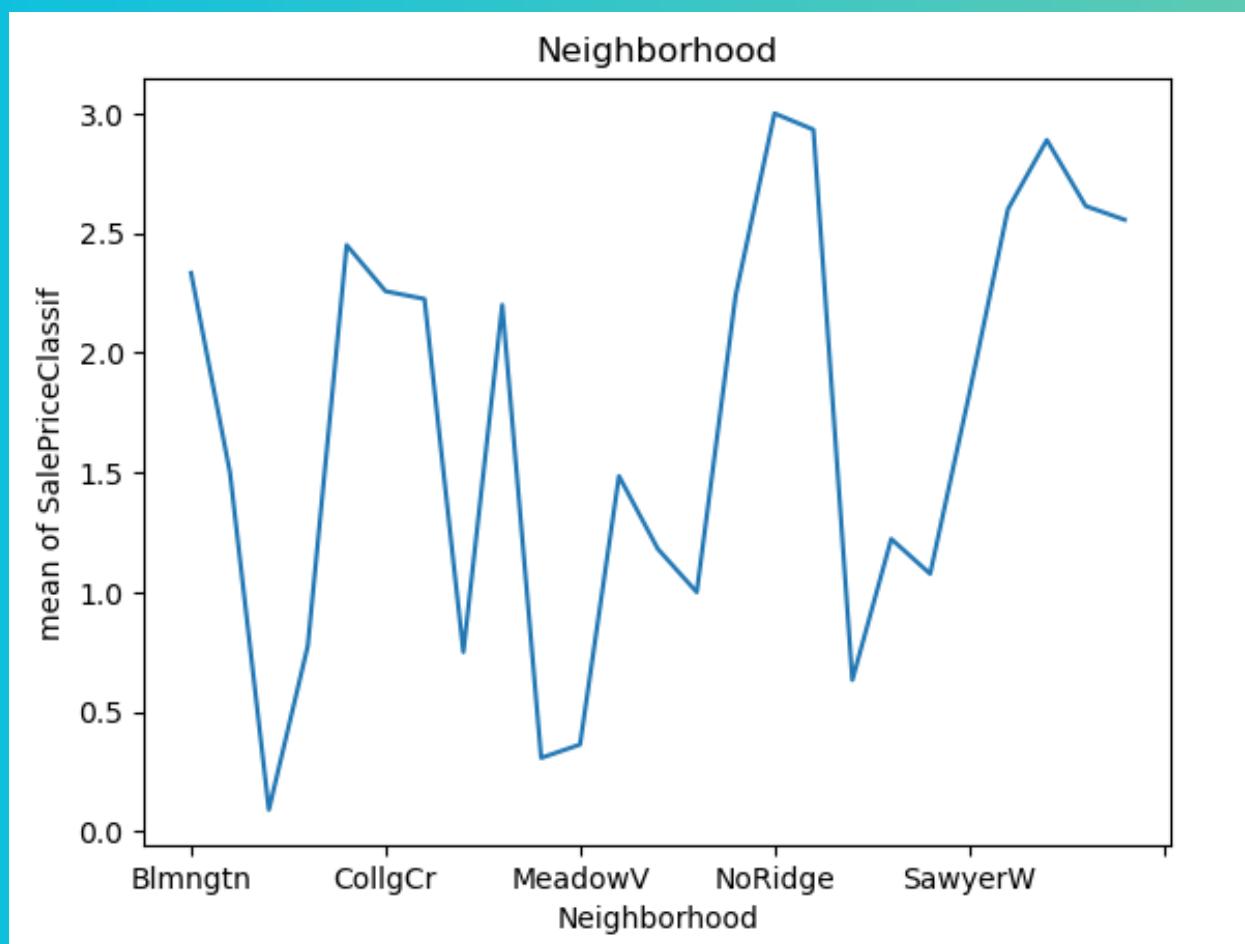
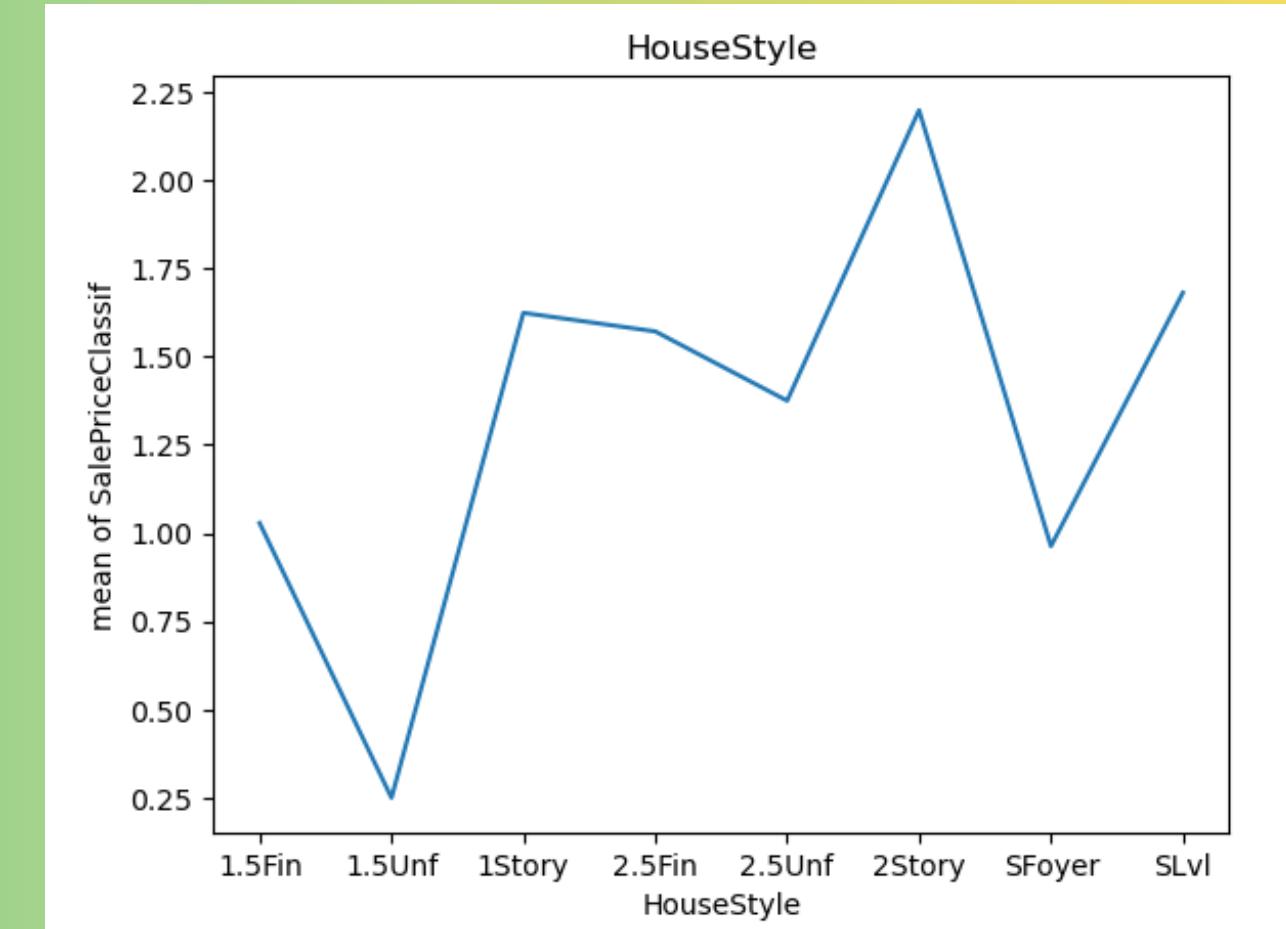
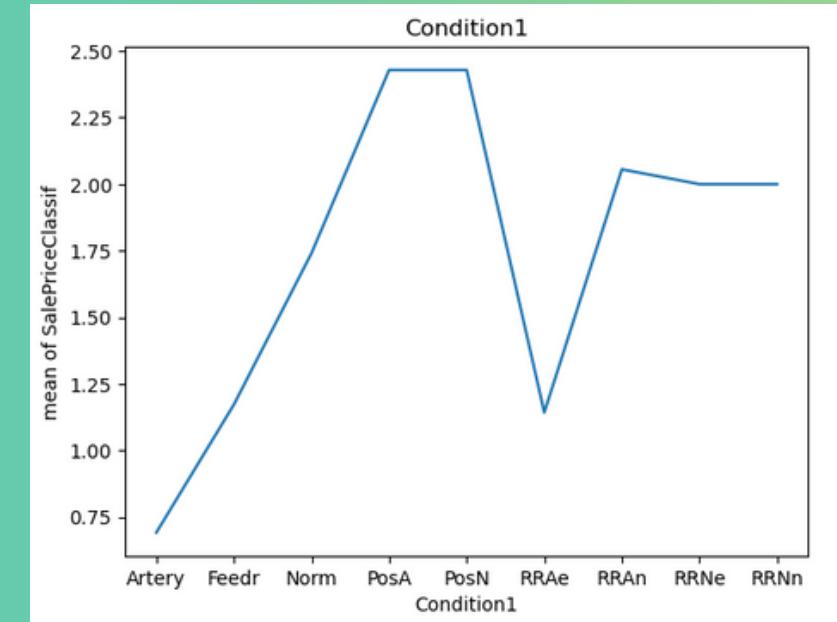
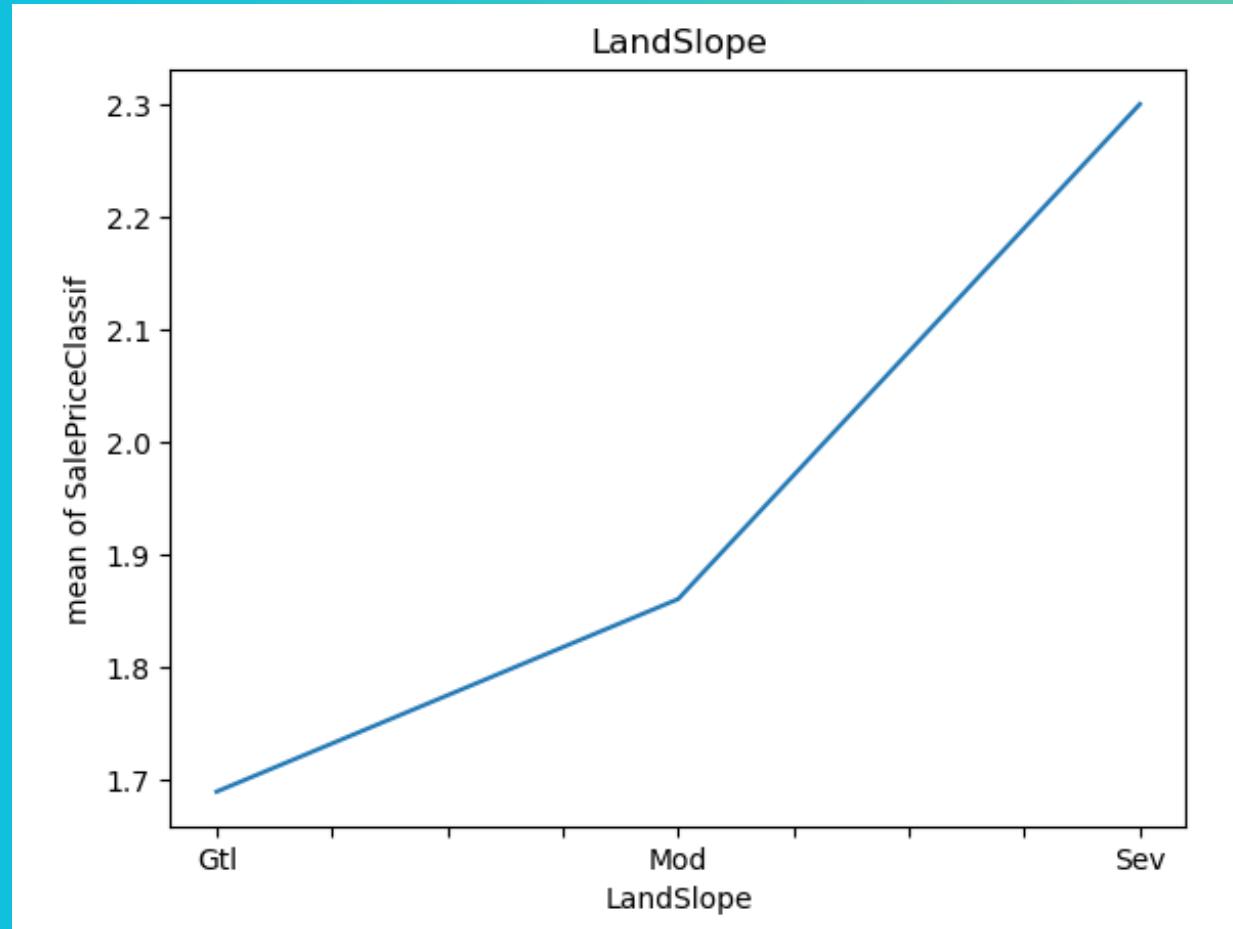
# Check Missing Value in Train and Test Sets

## 4.3 Encoding

### Checking relationships between categorical features and target

```
1 for var in categorical_features_train:  
2     check_relationship_with_target(X_train_impute, y_train, var)
```





```
1 categorical_norm_relationships = ['Street', 'Utilities', 'LandSlope', 'CentralAir']
2
3 categorical_not_norm_relationships = [
4     var for var in categorical_features_train if var not in categorical_norm_relationships
5 ]
6
7 print(' Categorical features with norm relationship with target: \n', categorical_norm_relationships, '\n')
8
9 print(' Categorical features with not norm relationship with target: \n ', categorical_not_norm_relationships, '\n')
```

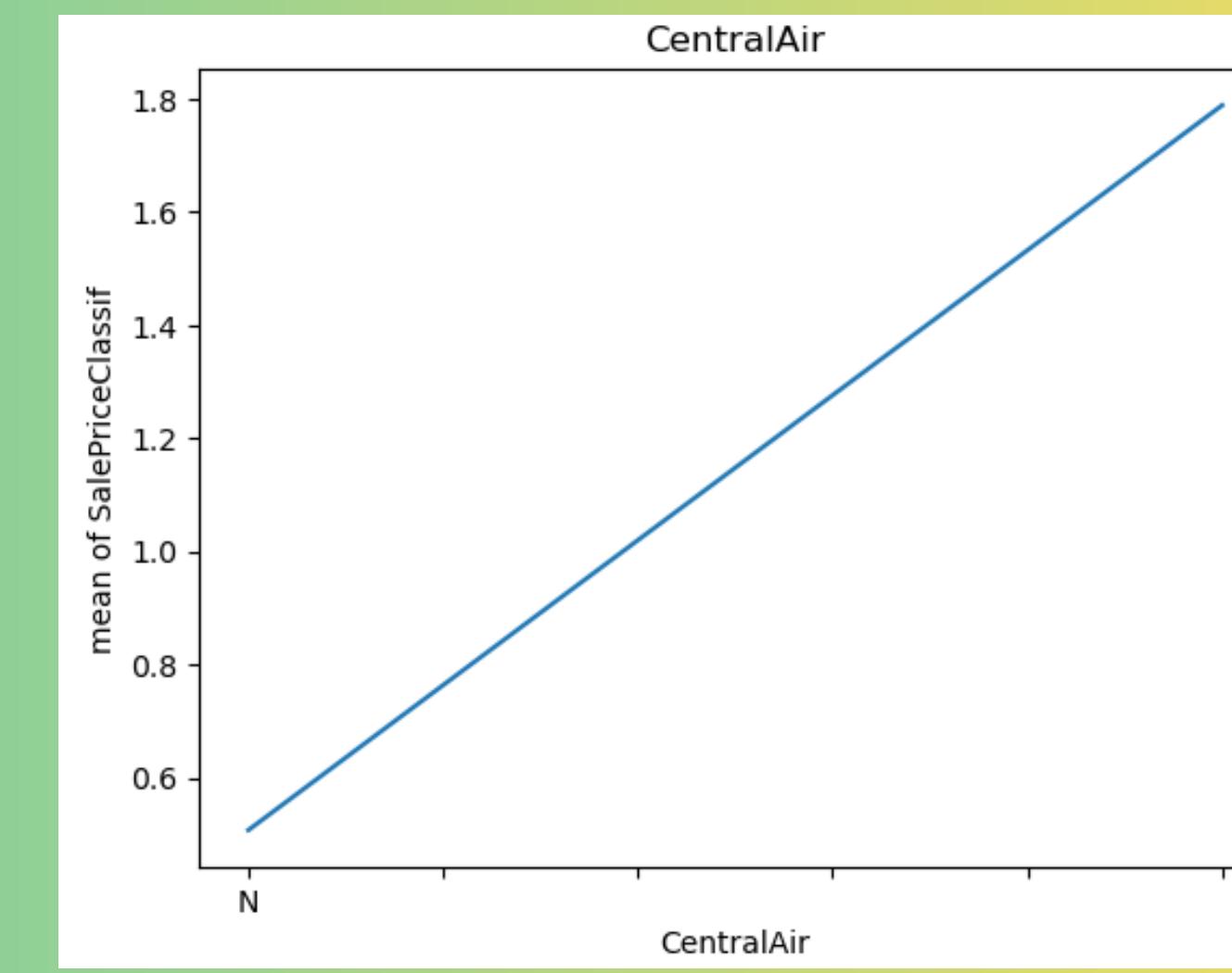
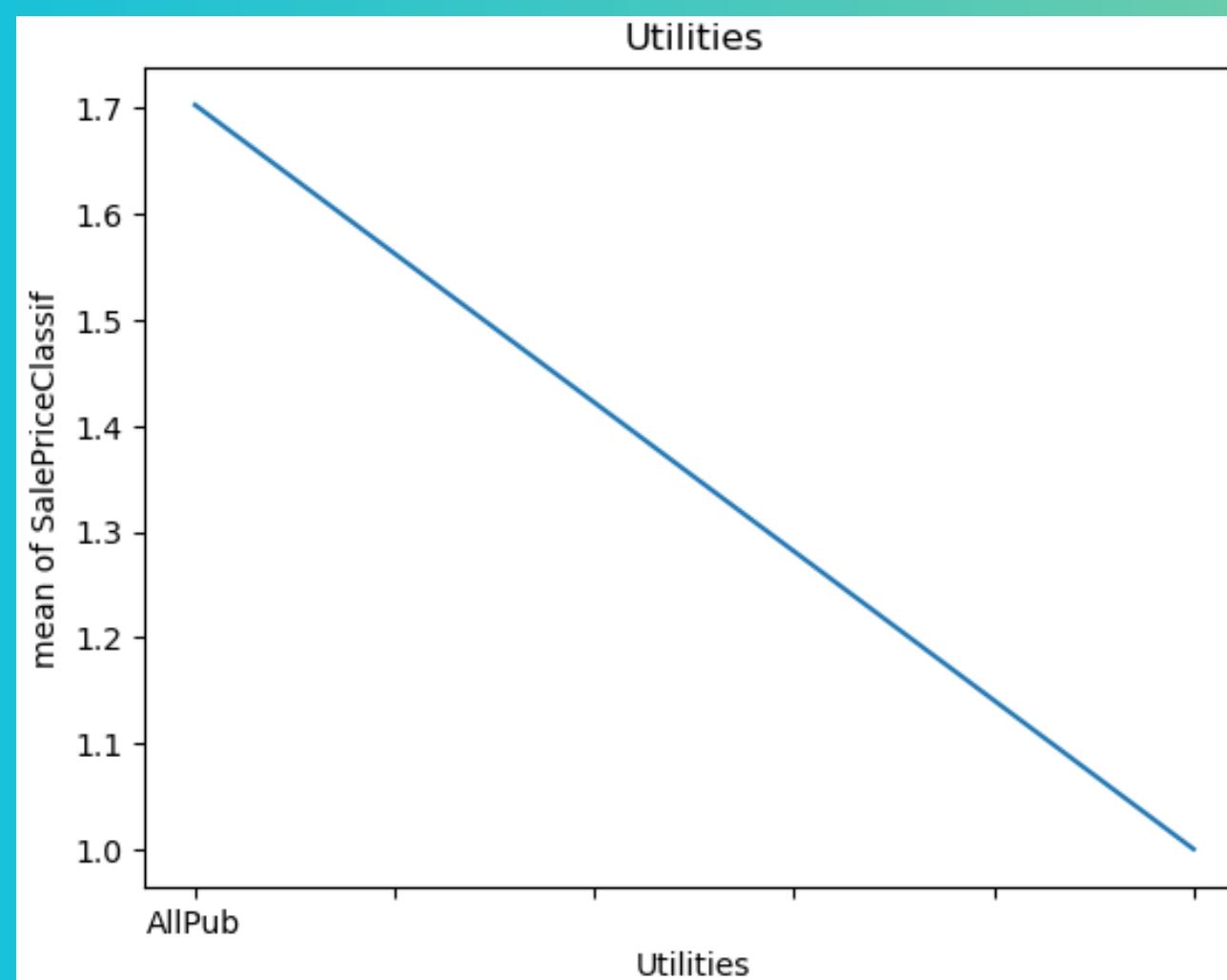
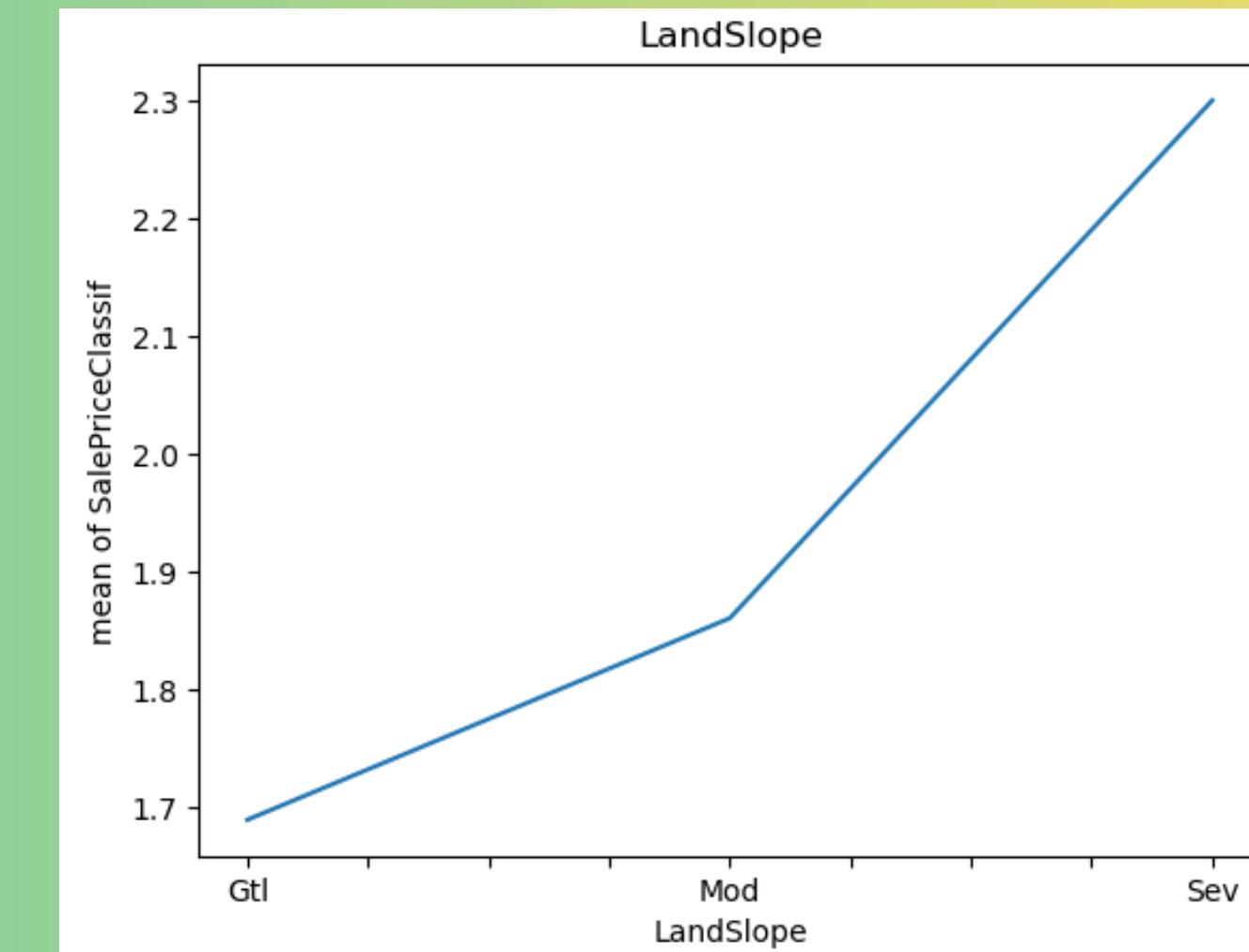
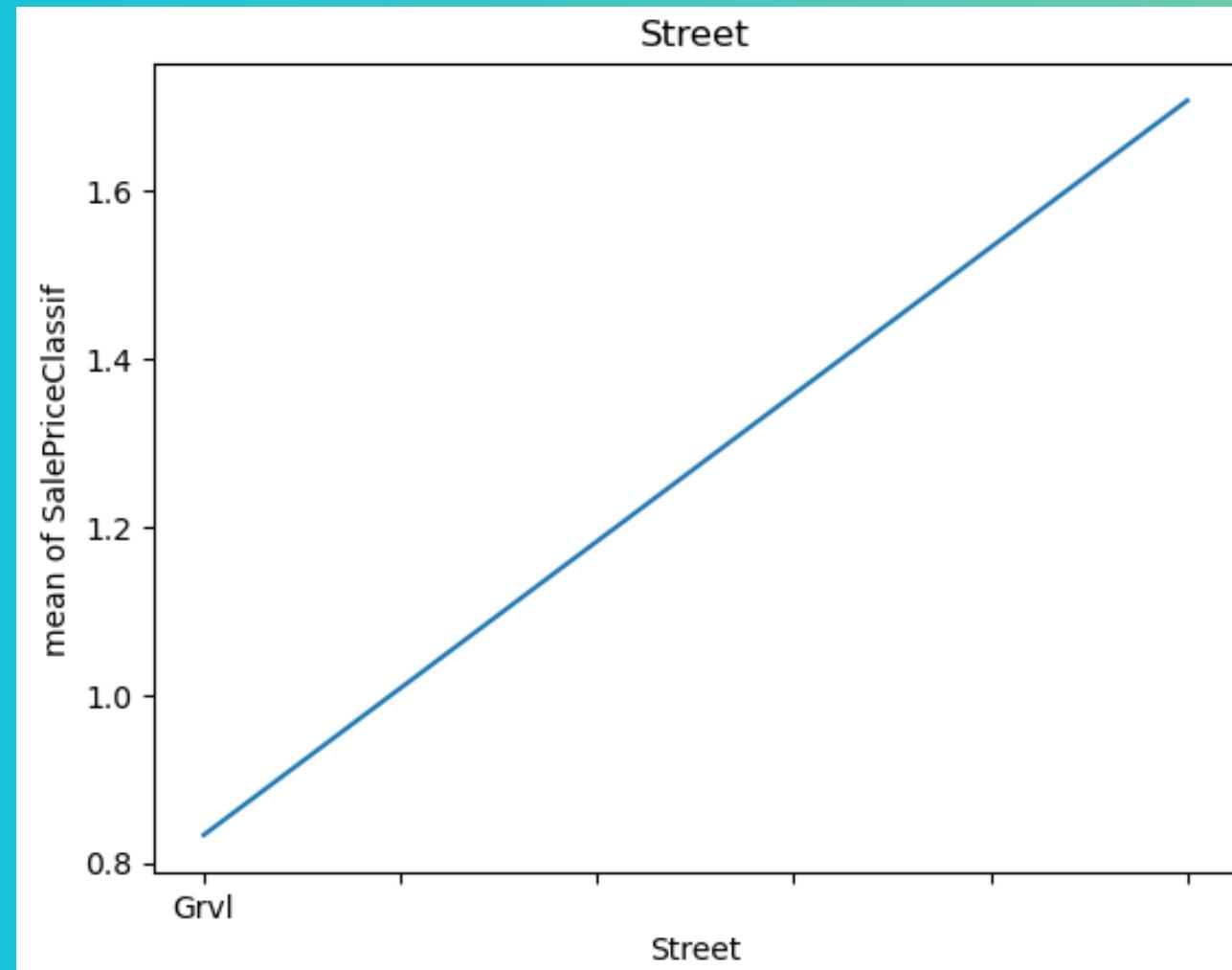
Categorical features with norm relationship with target:

```
['Street', 'Utilities', 'LandSlope', 'CentralAir']
```

Categorical features with not norm relationship with target:

```
['MSZoning', 'LotShape', 'LandContour', 'LotConfig', 'Neighborhood', 'Condition1', 'Condition2', 'BldgType', 'House
Style', 'RoofStyle', 'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'ExterQual', 'ExterCond', 'Foundation',
'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2', 'Heating', 'HeatingQC', 'Electrical',
'KitchenQual', 'Functional', 'FireplaceQu', 'GarageType', 'GarageFinish', 'GarageQual', 'GarageCond', 'PavedDrive',
'SaleType', 'SaleCondition']
```

# Categorical Features with Norm and Not Norm Relationships



## One hot encoding

```
1 # Cardinality for OneHot encoding
2
3 for var in categorical_norm_relationships:
4     print(var, '\n Count of unique values: ', X_train[var].nunique(), '\n')
```

Street

Count of unique values: 2

Utilities

Count of unique values: 2

LandSlope

Count of unique values: 3

CentralAir

Count of unique values: 2

```
1 # Don't tune top_categories, for the reason of low cardinality
2 enc_ohe = OneHotEncoder(
3     variables = categorical_norm_relationships,
4     drop_last = False)
5
6 enc_ohe.fit(X_train_impute)
7
8 X_train_enc_ohe = enc_ohe.transform(X_train_impute)
9
10 X_test_enc_ohe = enc_ohe.transform(X_test_impute)
```

```
1 len(X_train_impute.columns), len(X_train_enc_ohe.columns)
```

(74, 79)

## Mean encoding

```
1 for var in X_train_impute.columns:  
2     if X_train_impute[var].isnull().mean() > 0:  
3         print(var)  
  
1 categorical_mean = categorical_not_norm_relationships.copy()  
2  
3 categorical_mean.remove('Condition1')  
4 categorical_mean.remove('Condition2')  
5 categorical_mean.remove('RoofMatl')  
6 categorical_mean.remove('Exterior1st')  
7 categorical_mean.remove('Exterior2nd')  
8 categorical_mean.remove('ExterCond')  
9 categorical_mean.remove('HeatingQC')  
10 categorical_mean.remove('Functional')  
11 categorical_mean.remove('Neighborhood')  
12 categorical_mean.remove('Electrical')  
13 categorical_mean.remove('RoofStyle')  
14 categorical_mean.remove('Heating')  
15 categorical_mean.remove('GarageQual')  
16  
17 encoder_mean = MeanEncoder(variables = categorical_mean)  
18  
19 encoder_mean.fit(X_train_impute, y_train)  
20  
21 X_train_enc_mean = encoder_mean.transform(X_train_impute)  
22  
23 X_test_enc_mean = encoder_mean.transform(X_test_impute)
```

We remove some features because encoding by mean value creates empty variables in the test set

```
1 categorical_problem = [  
2     var for var in categorical_not_norm_relationships if var not in categorical_mean  
3 ]  
4  
5 categorical_problem  
  
['Neighborhood',  
 'Condition1',  
 'Condition2',  
 'RoofStyle',  
 'RoofMatl',  
 'Exterior1st',  
 'Exterior2nd',  
 'ExterCond',  
 'Heating',  
 'HeatingQC',  
 'Electrical',  
 'Functional',  
 'GarageQual']
```

## Categorical\_Problem

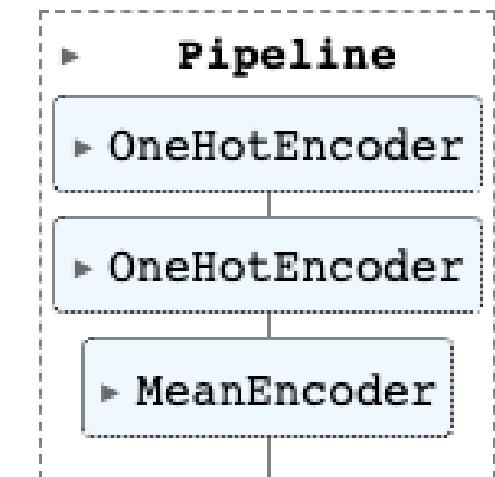
```
1 X_train_impute[categorical_problem].nunique()
```

Neighborhood	25
Condition1	9
Condition2	4
RoofStyle	6
RoofMatl	7
Exterior1st	14
Exterior2nd	15
ExterCond	4
Heating	5
HeatingQC	5
Electrical	5
Functional	7
GarageQual	6
<b>dtype:</b>	int64

## Categorical\_Problem nunique

# Build Pipe Line Encoding

```
1 for var in categorical_features_train:  
2     if X_train_impute[var].dtype != 'O':  
3         print(var)  
  
1 pipeline_encoding = Pipeline(  
2     steps = [  
3         ('ohe', OneHotEncoder(variables = categorical_norm_relationships)),  
4         ('ohe_top_cat', OneHotEncoder(variables = categorical_problem, top_categories = 3)),  
5         ('mean', MeanEncoder(variables = categorical_mean))  
6     ]  
7  
8 pipeline_encoding.fit(X_train_impute, y_train)
```



```
1 X_train_impute_enc = pipeline_encoding.transform(X_train_impute)  
2  
3 X_test_impute_enc = pipeline_encoding.transform(X_test_impute)
```

## Checking object features in set

```
1 # Checking object features in sets  
2  
3 for var in X_train_impute_enc.columns:  
4     if X_train_impute_enc[var].dtype == 'O':  
5         print(var)
```

# Discretization

## Discretization

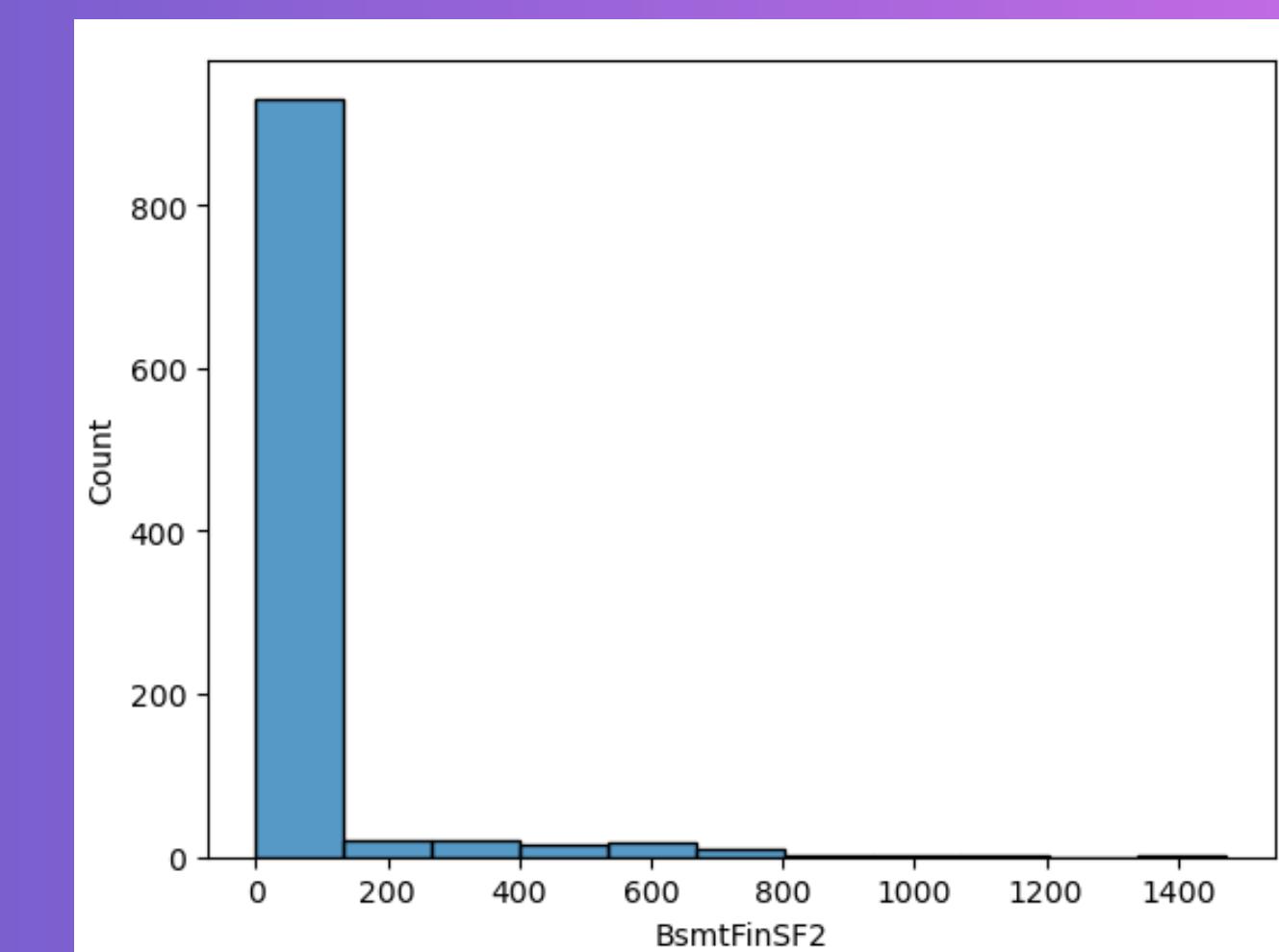
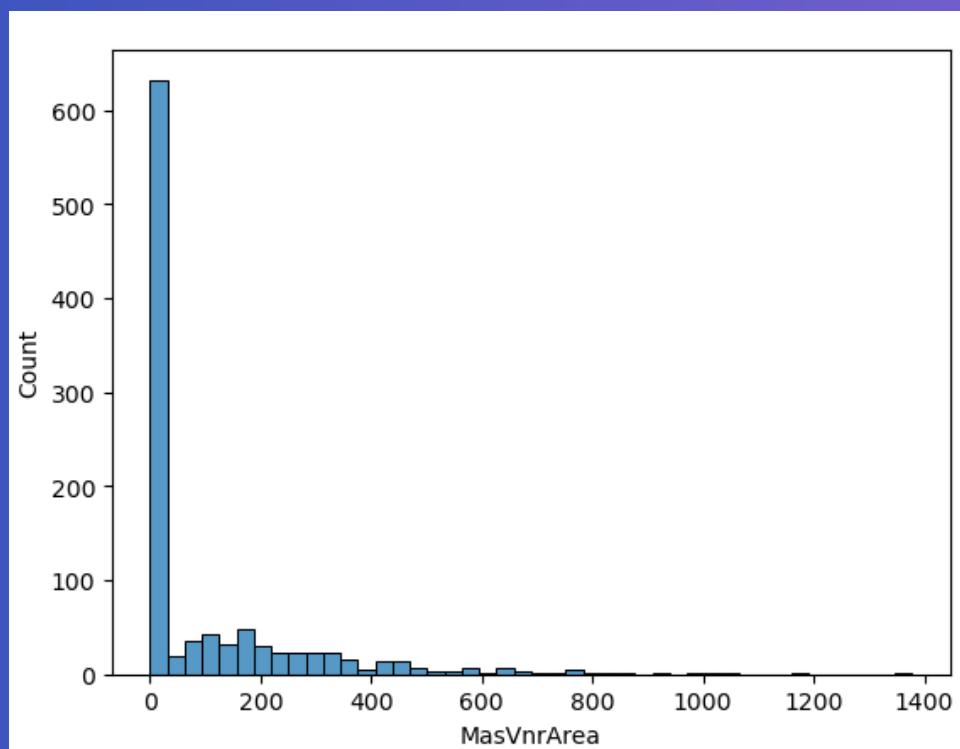
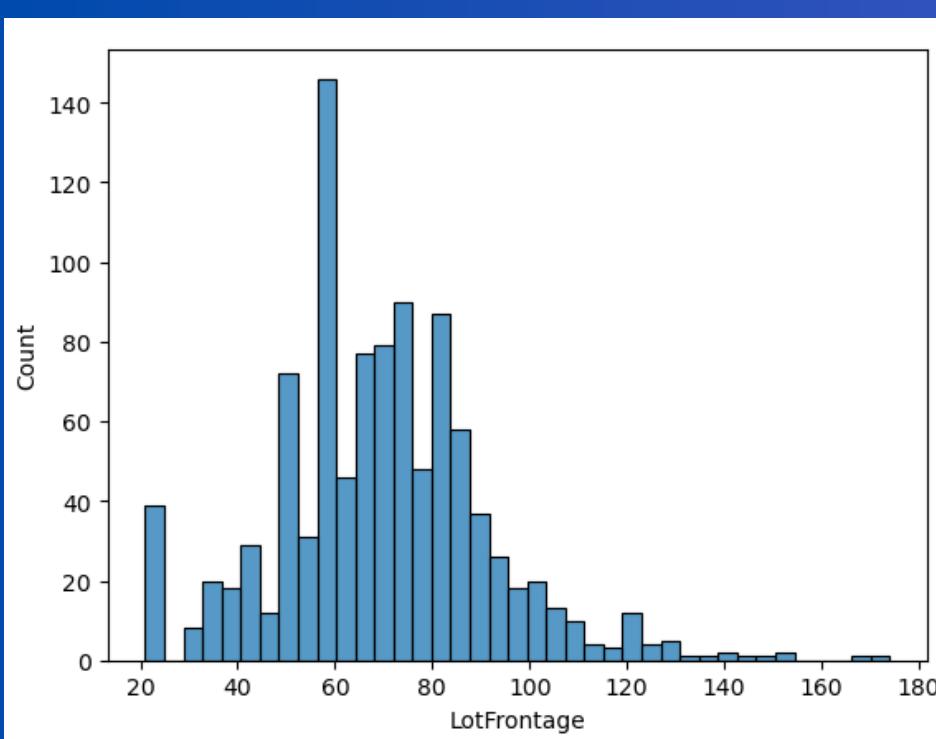
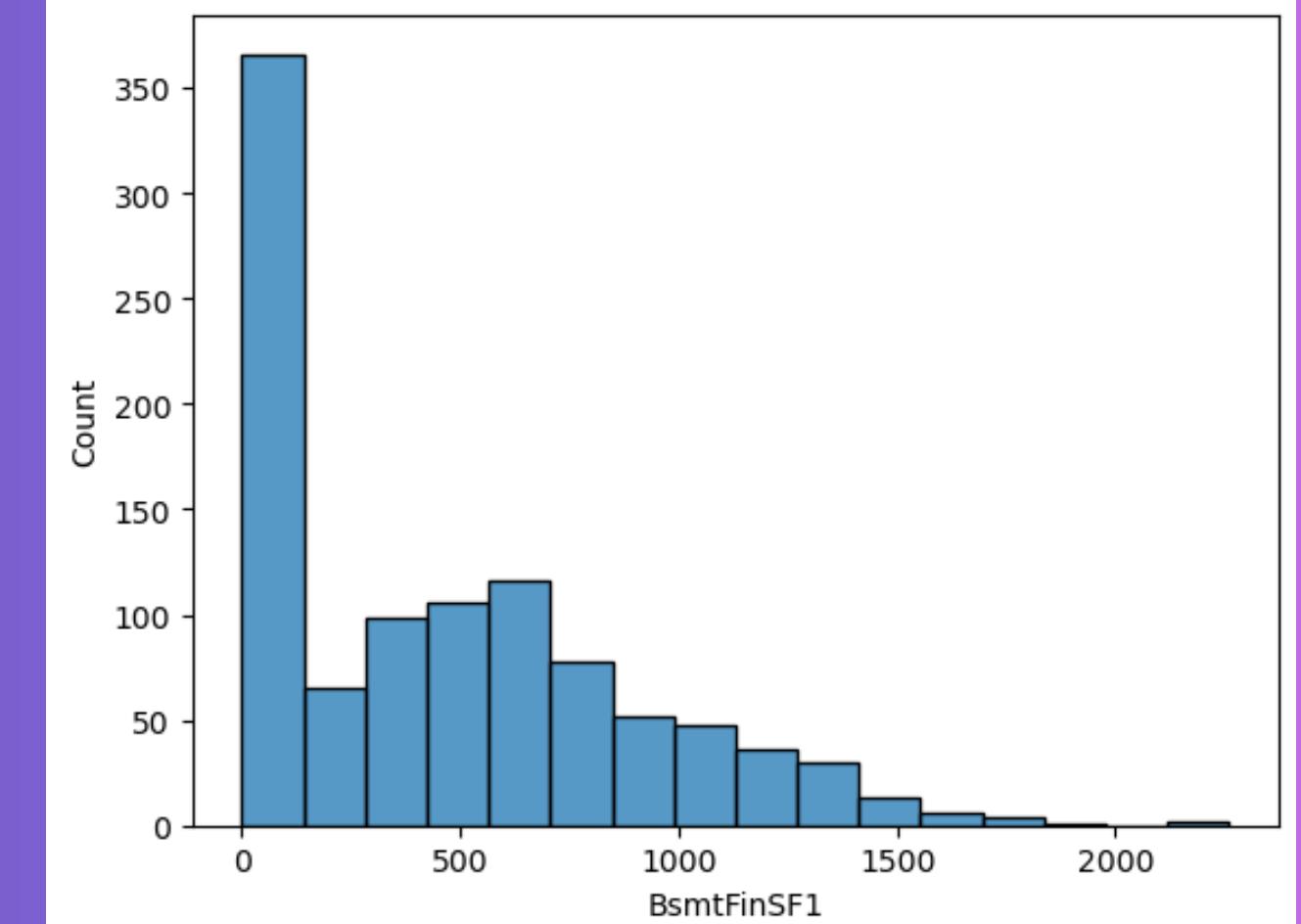
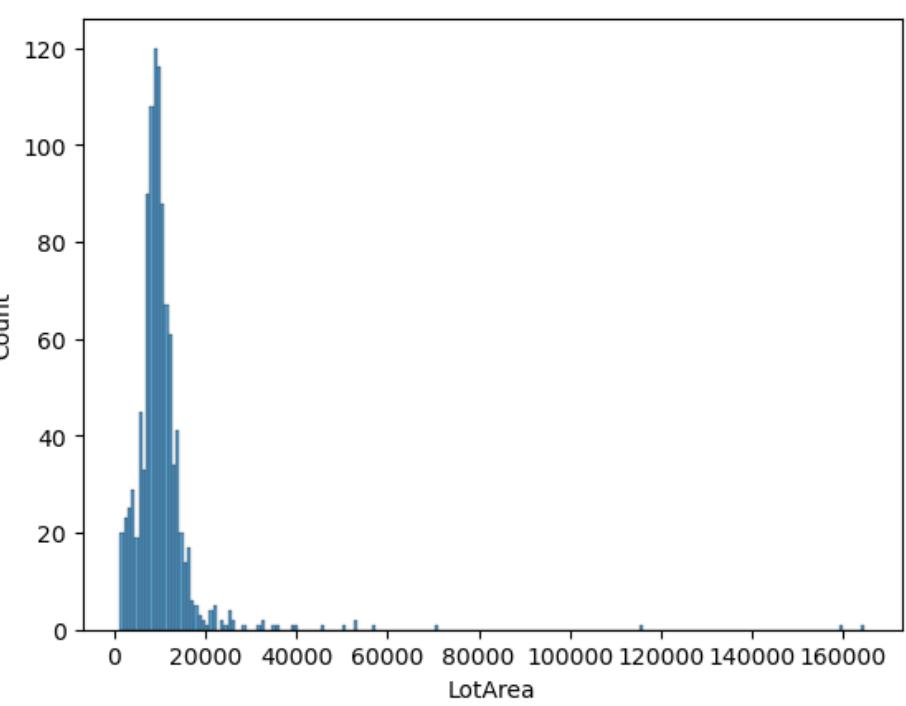
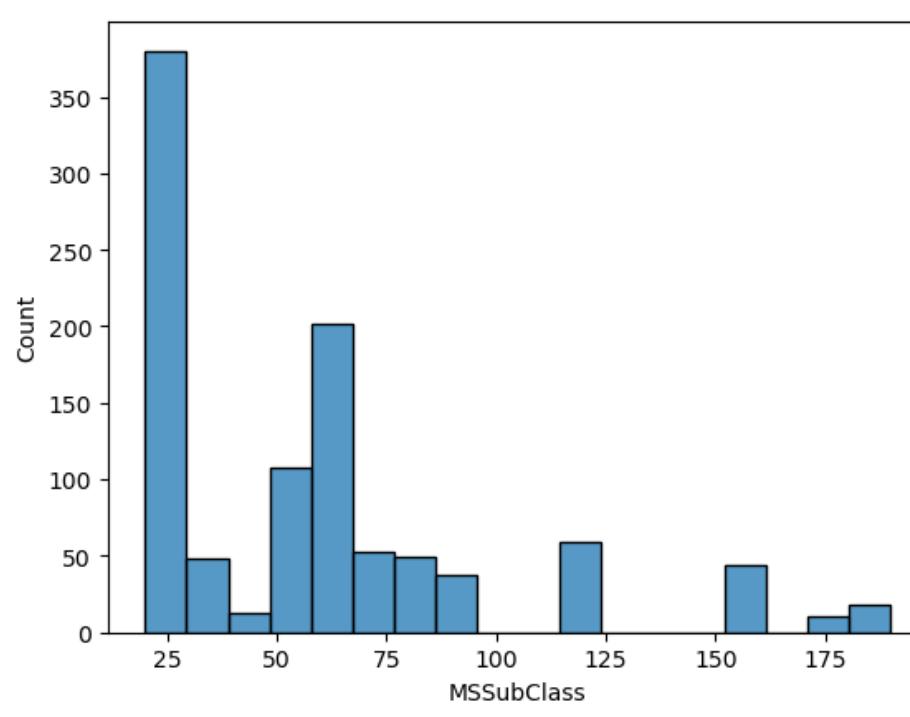
```
1 # Make continuous features list
2
3 continuous_features_train = [
4     var for var in X_train_impute_enc.columns if X_train_impute_enc[var].nunique() > 10
5 ]
6
7 continuous_features_train
```

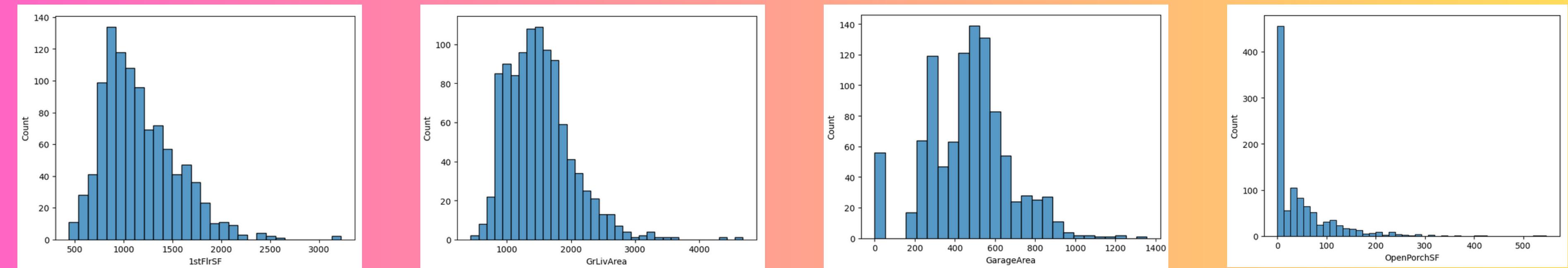
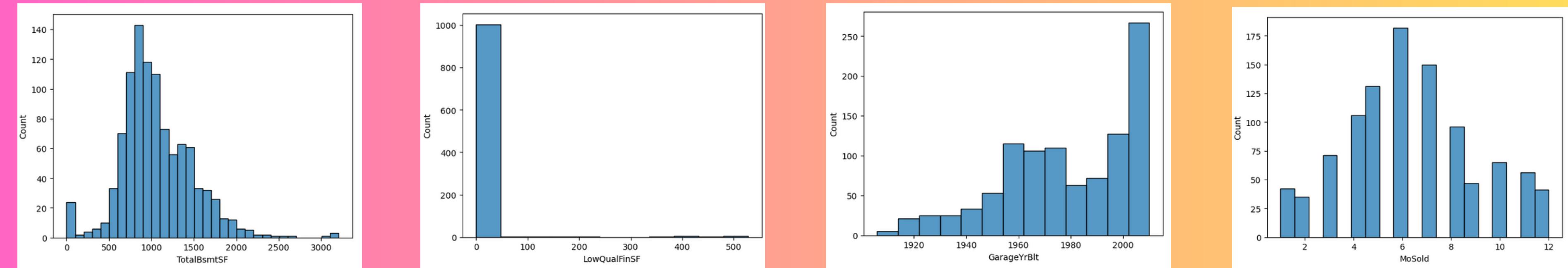
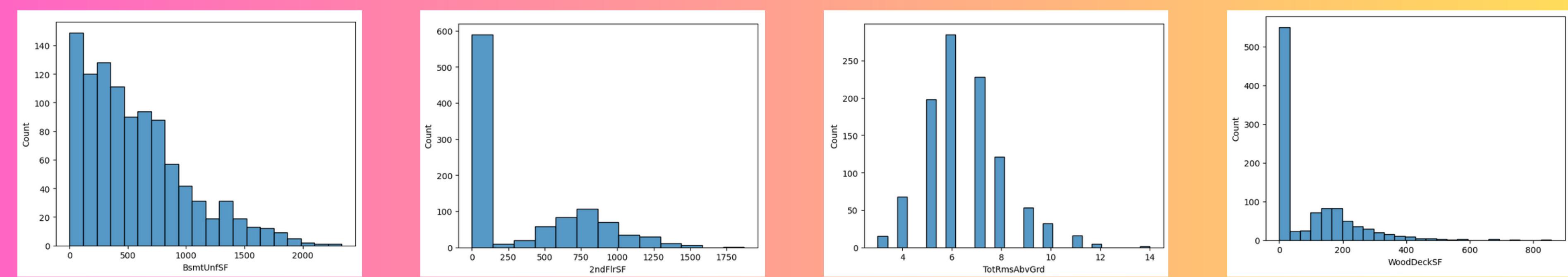
```
['MSSubClass',
 'LotFrontage',
 'LotArea',
 'MasVnrArea',
 'BsmtFinSF1',
 'BsmtFinSF2',
 'BsmtUnfSF',
 'TotalBsmtSF',
 '1stFlrSF',
 '2ndFlrSF',
 'LowQualFinSF',
 'GrLivArea',
 'TotRmsAbvGrd',
 'GarageYrBlt',
 'GarageArea',
 'WoodDeckSF',
 'OpenPorchSF',
 'EnclosedPorch',
 '3SsnPorch',
 'ScreenPorch',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

## Continuous Features List

```
# Check distribution of continuous features
```

```
for var in continuous_features_train:  
    sns.histplot(X_train_impute_enc[var])  
    plt.show()
```





# Equal Width Discretization

```
equal_width = EqualWidthDiscretiser(bins = 4, variables = continuous_features_train)

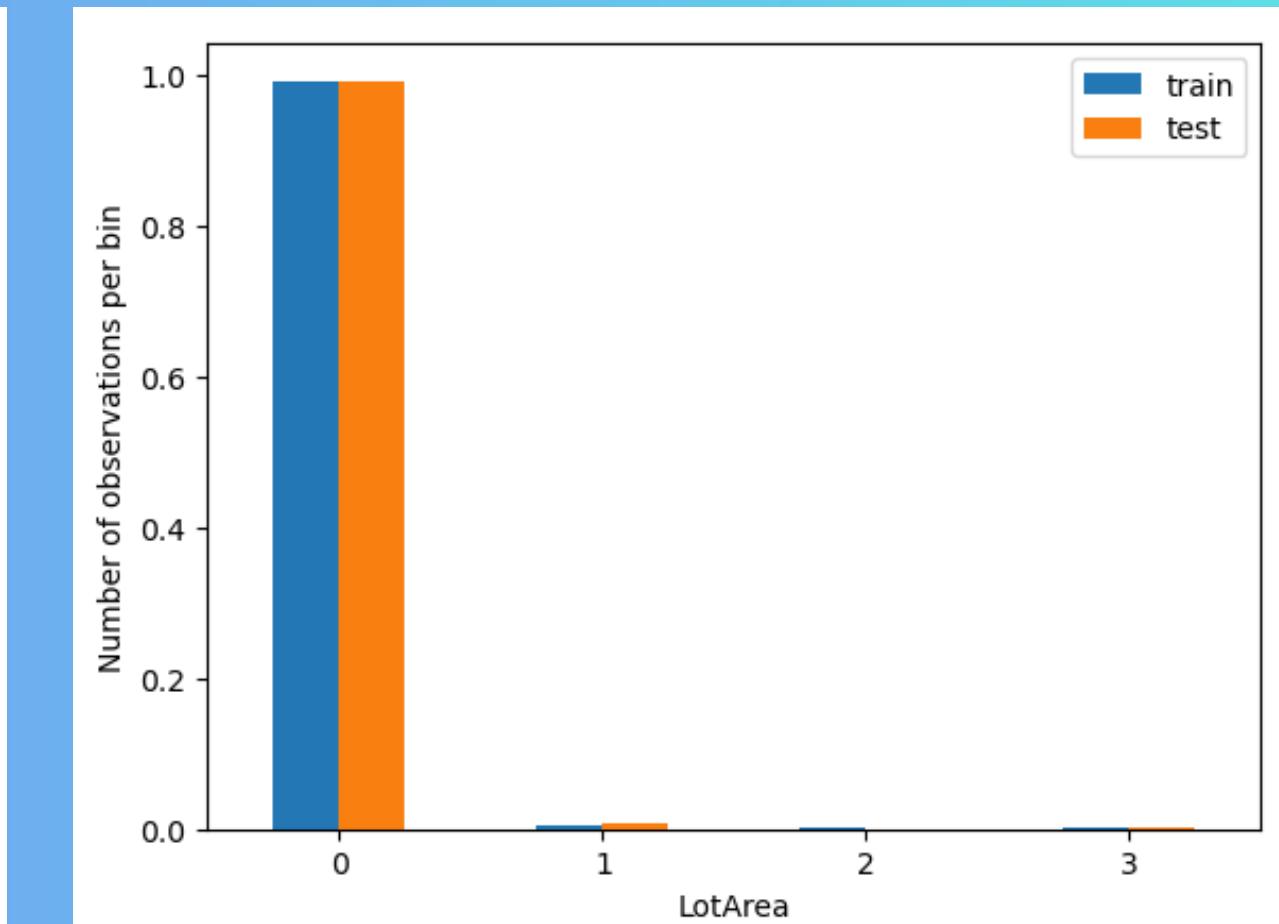
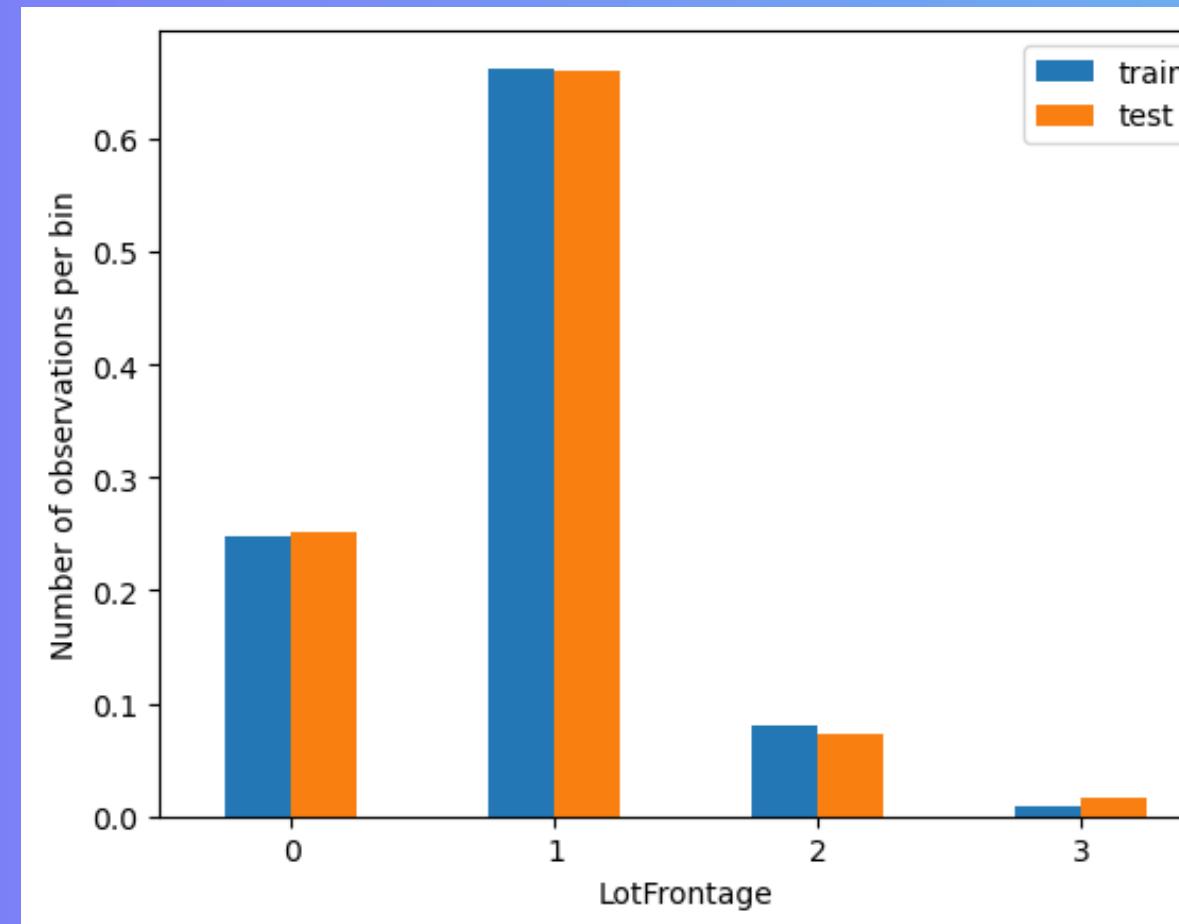
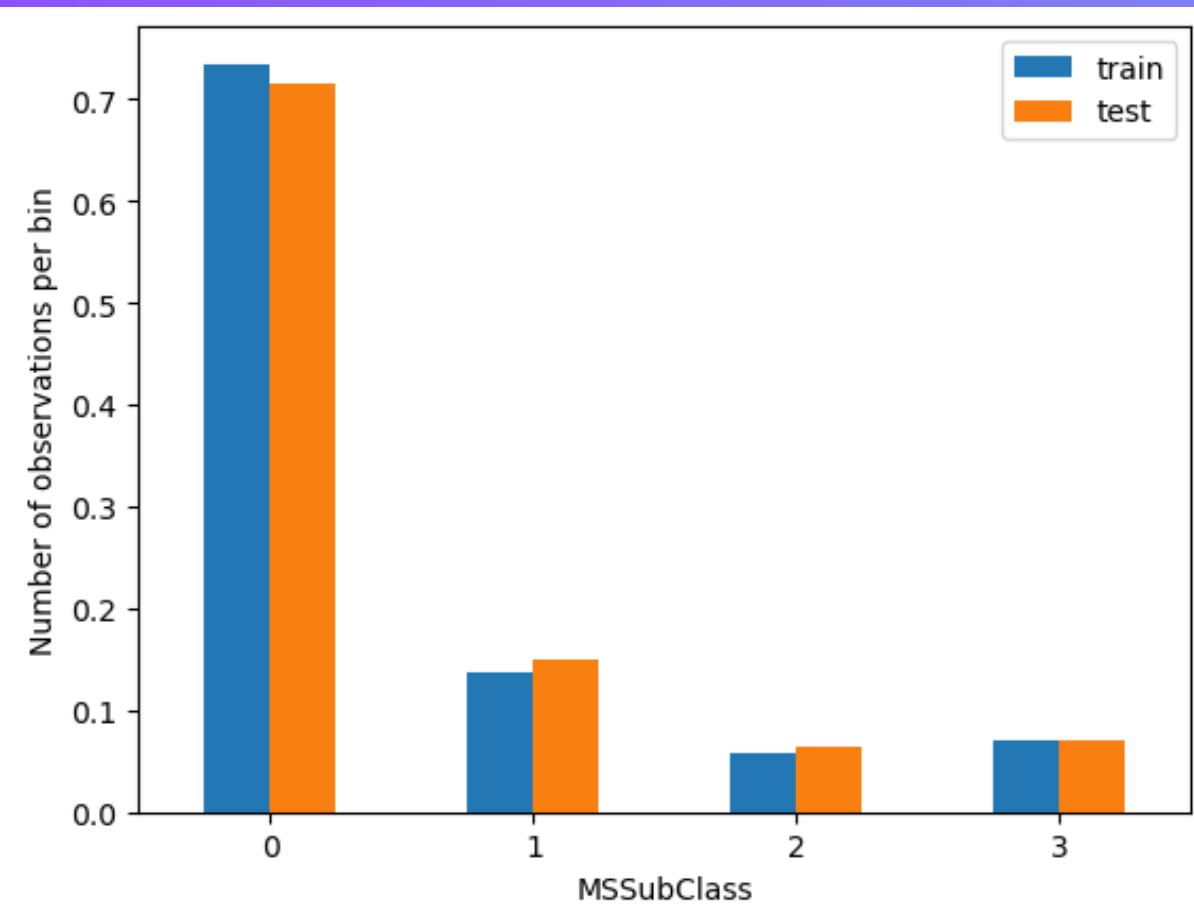
equal_width.fit(X_train_impute_enc)

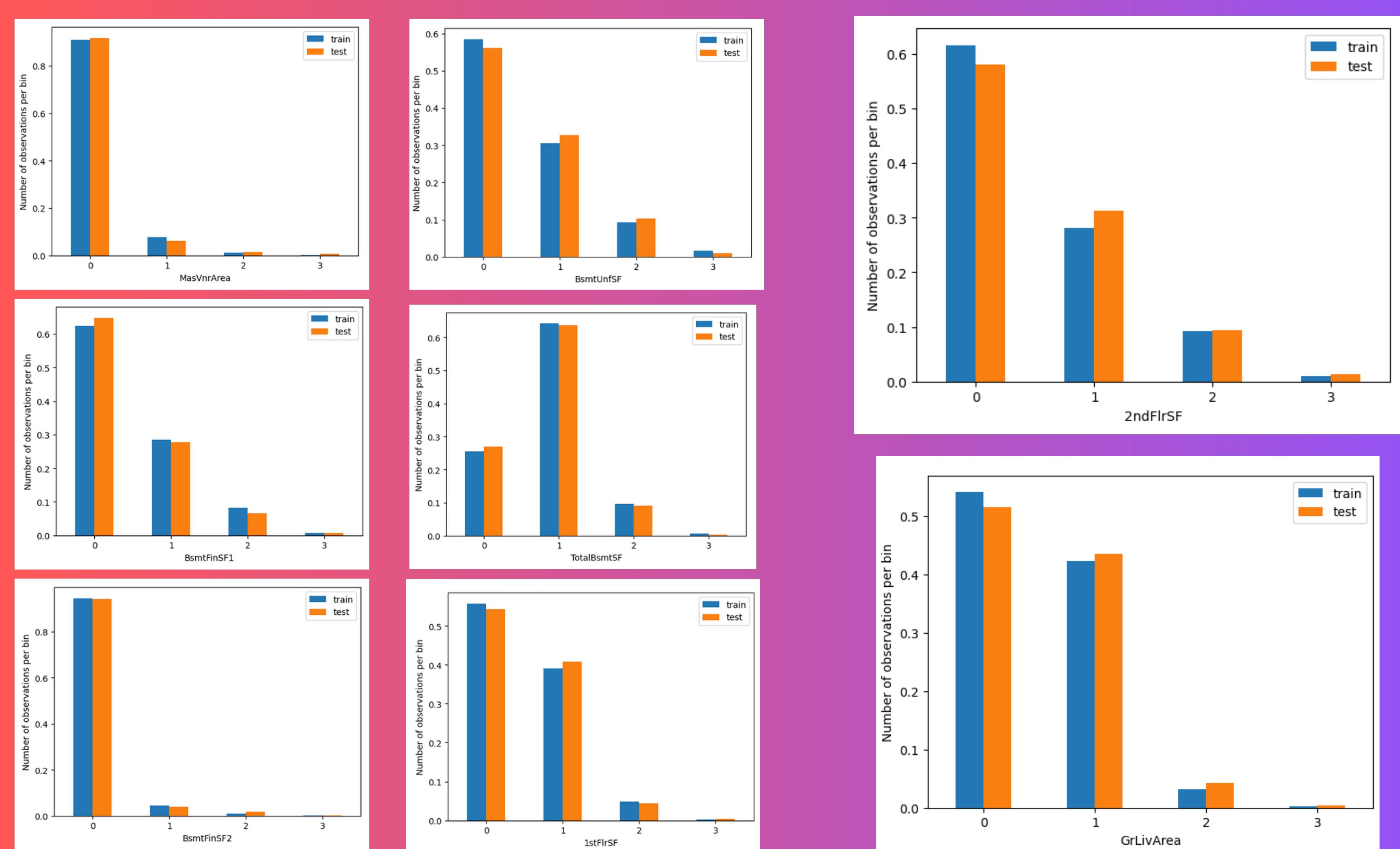
X_train_equal_width = equal_width.transform(X_train_impute_enc)

X_test_equal_width = equal_width.transform(X_test_impute_enc)
```

## Check balancing of data between train and test sets

```
for var in continuous_features_train:
    check_balance_of_discrete_and_categorical(X_train_equal_width, X_test_equal_width, var)
```



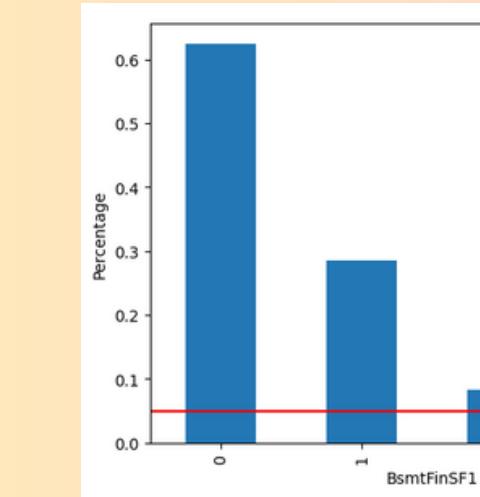
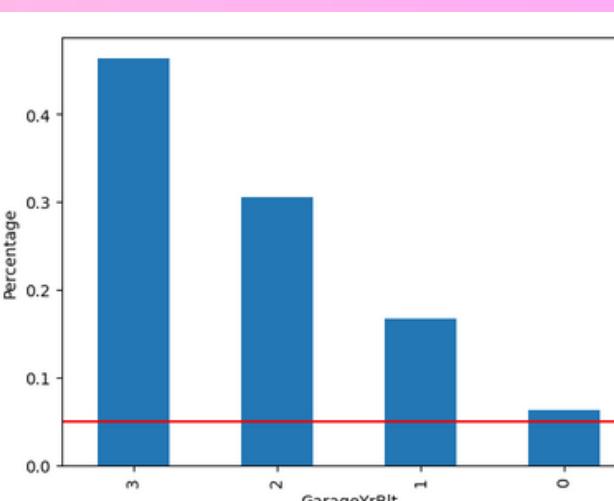
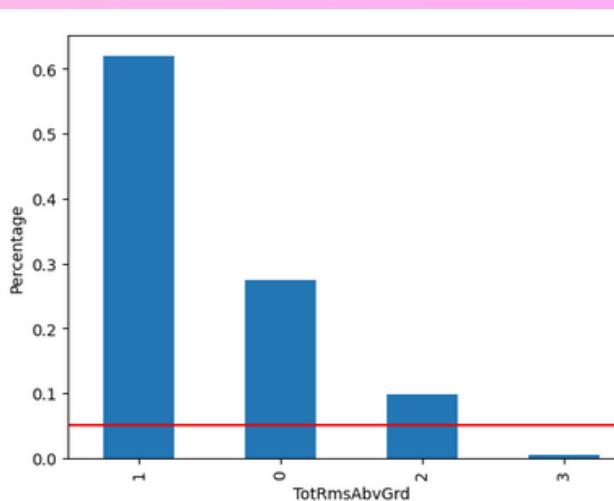
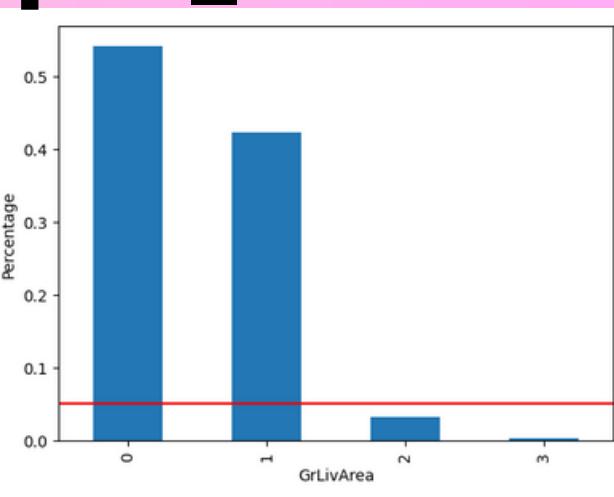
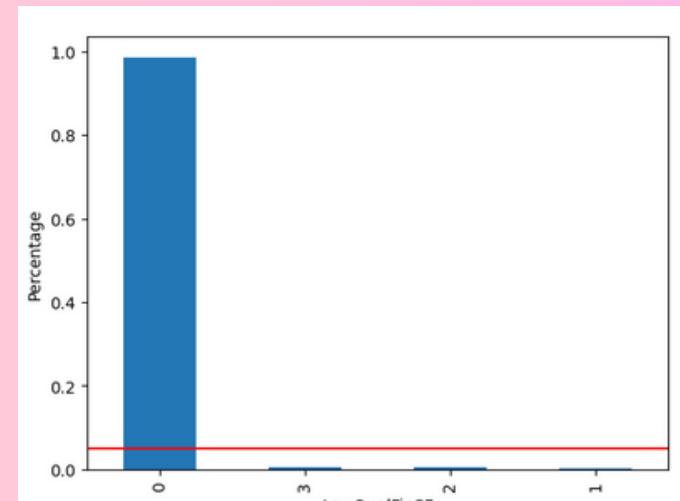
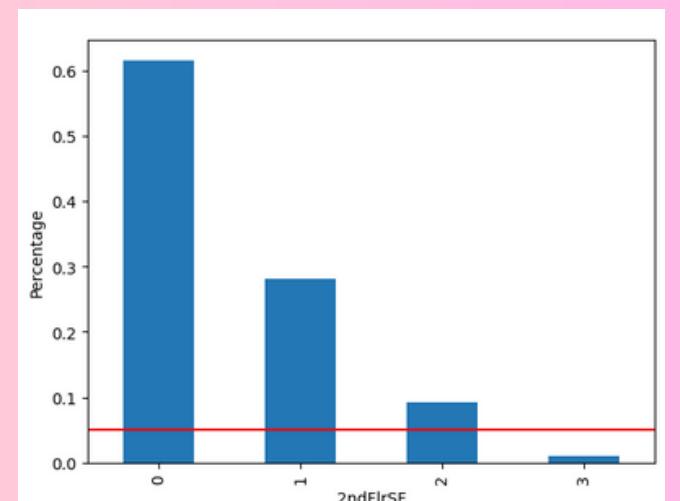
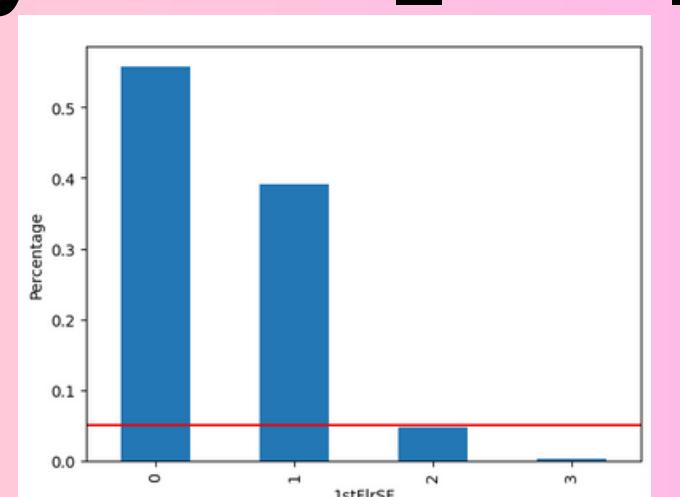
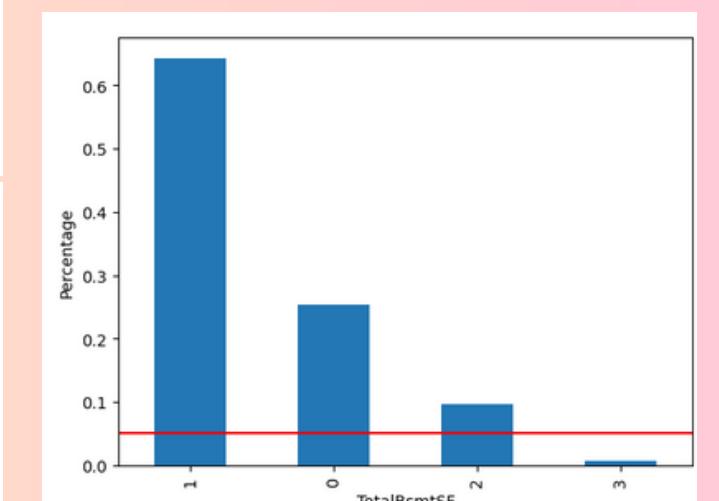
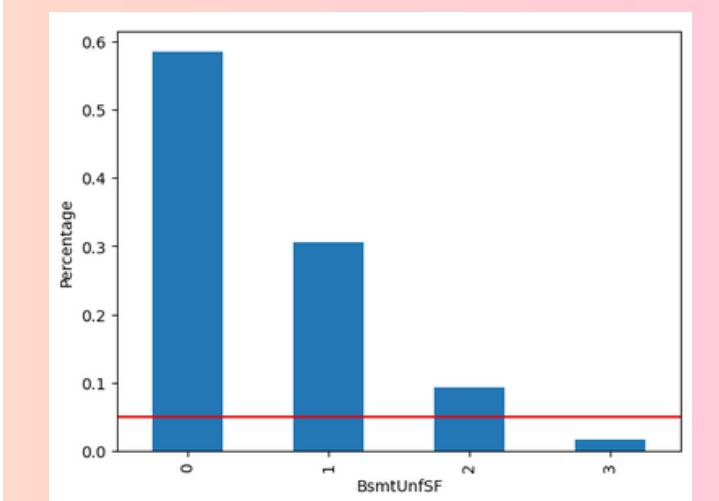
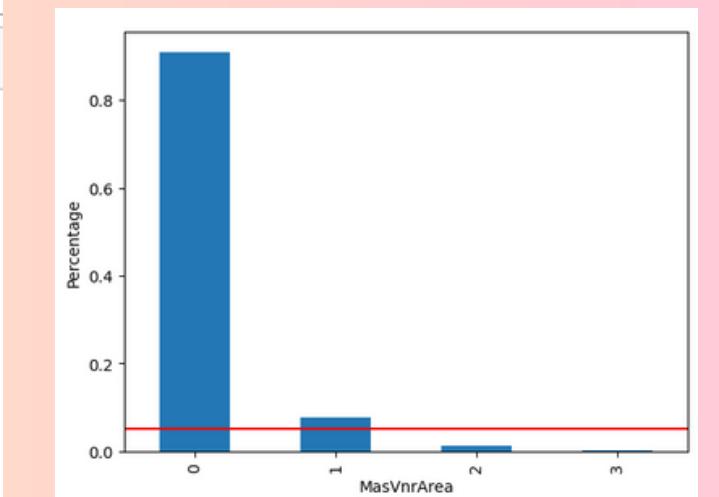
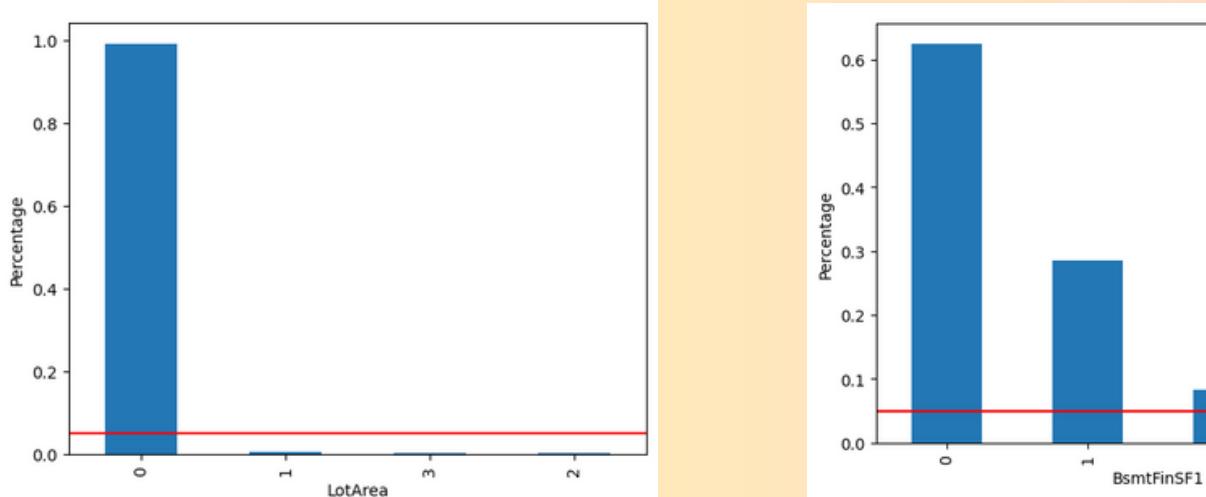
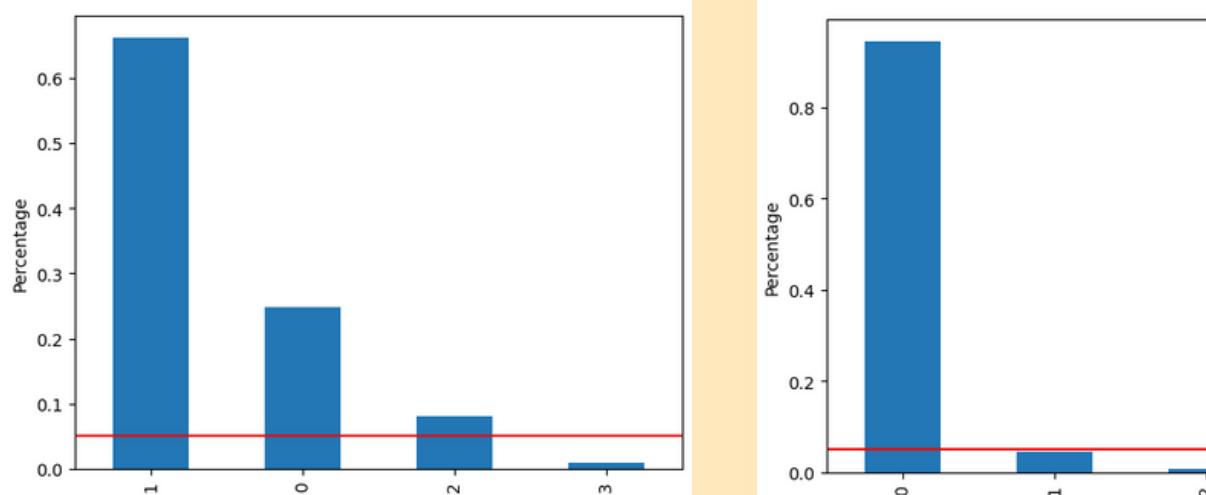
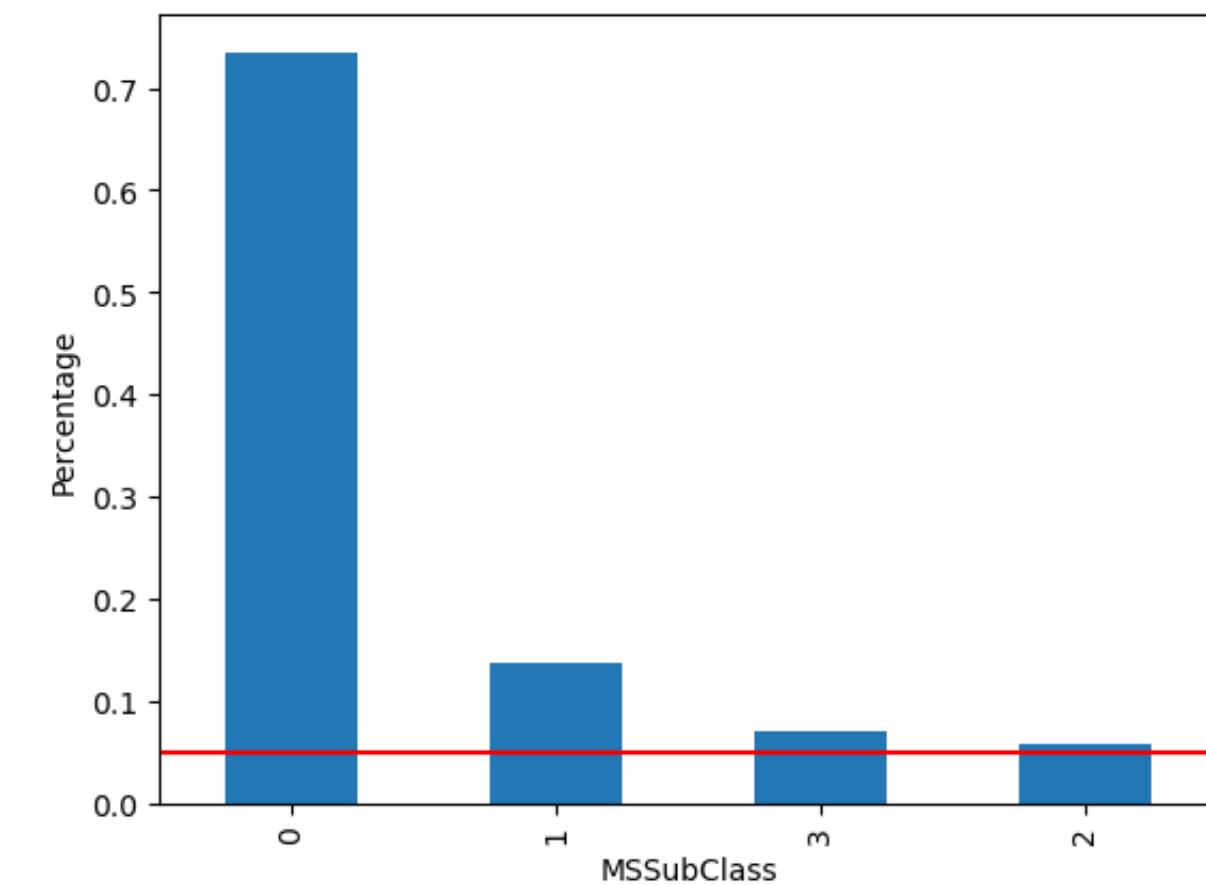


### Rare labels analysis

# Rare Labels Analysis

## X\_train\_equal\_width

```
1 check_rare_labels_analysis(X_train_equal_width, continuous_features_train)
```

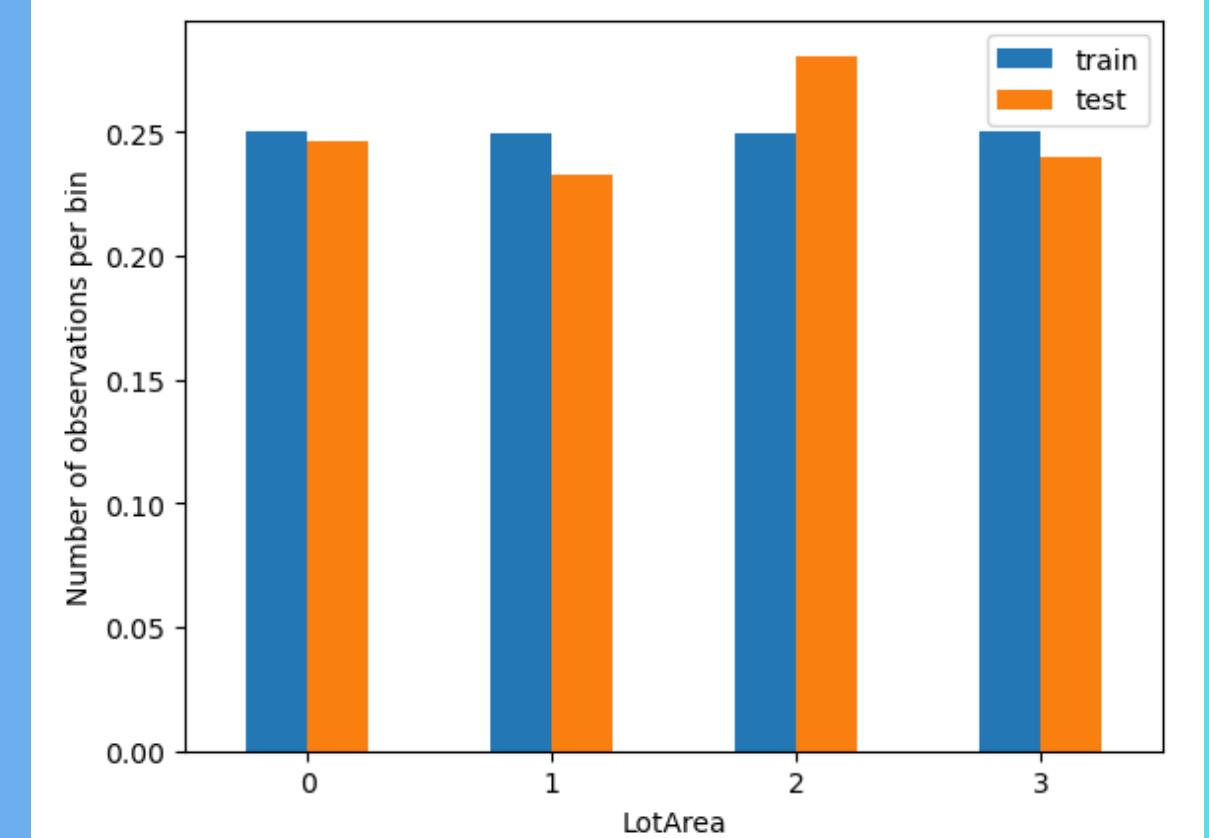
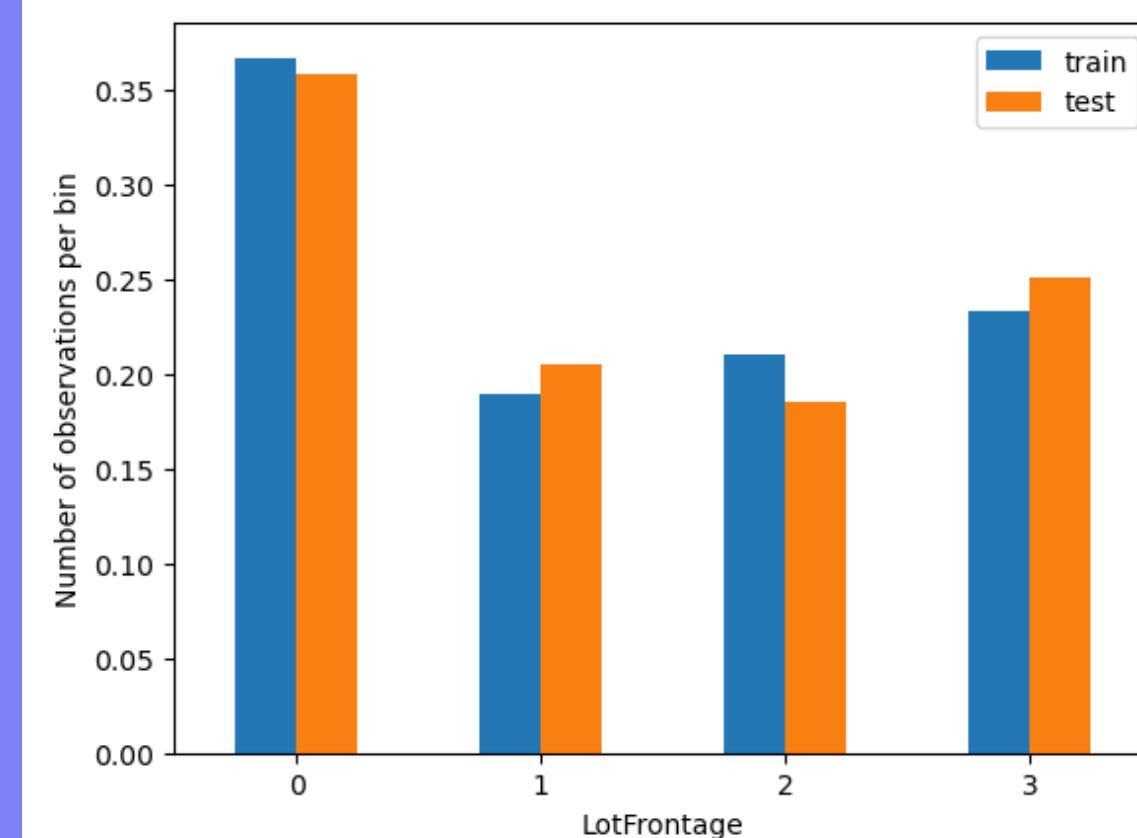
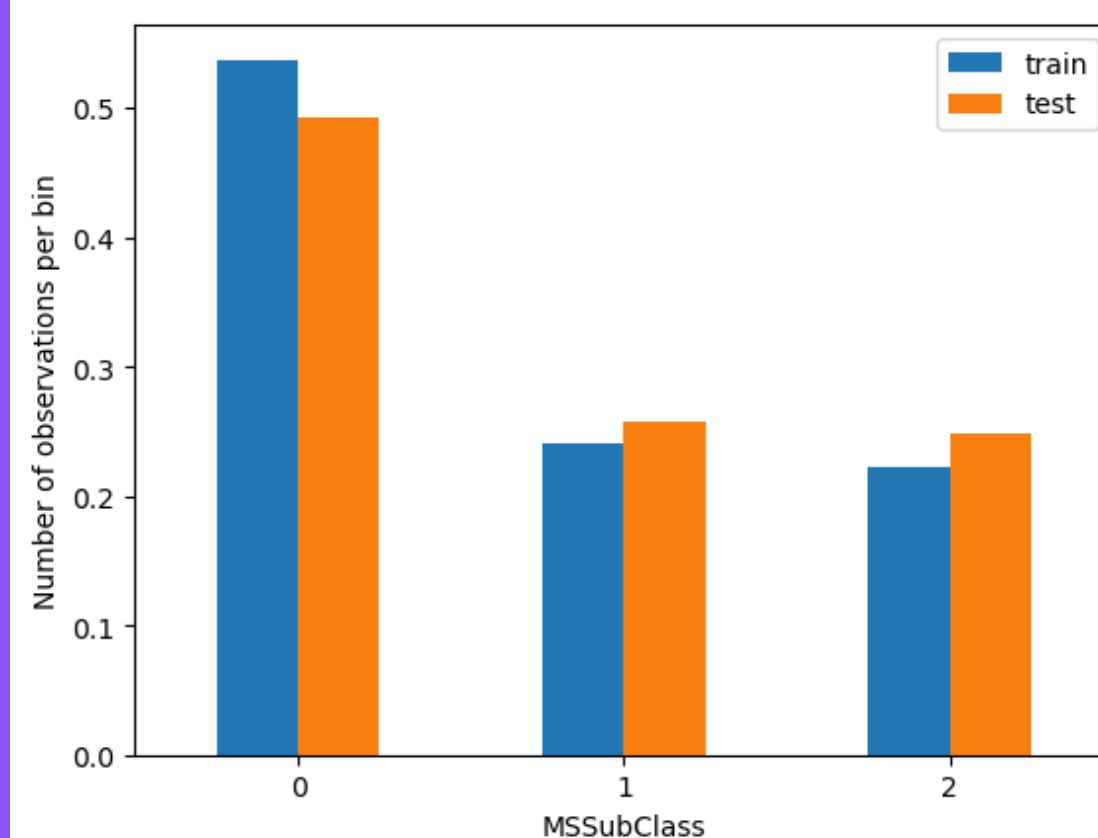


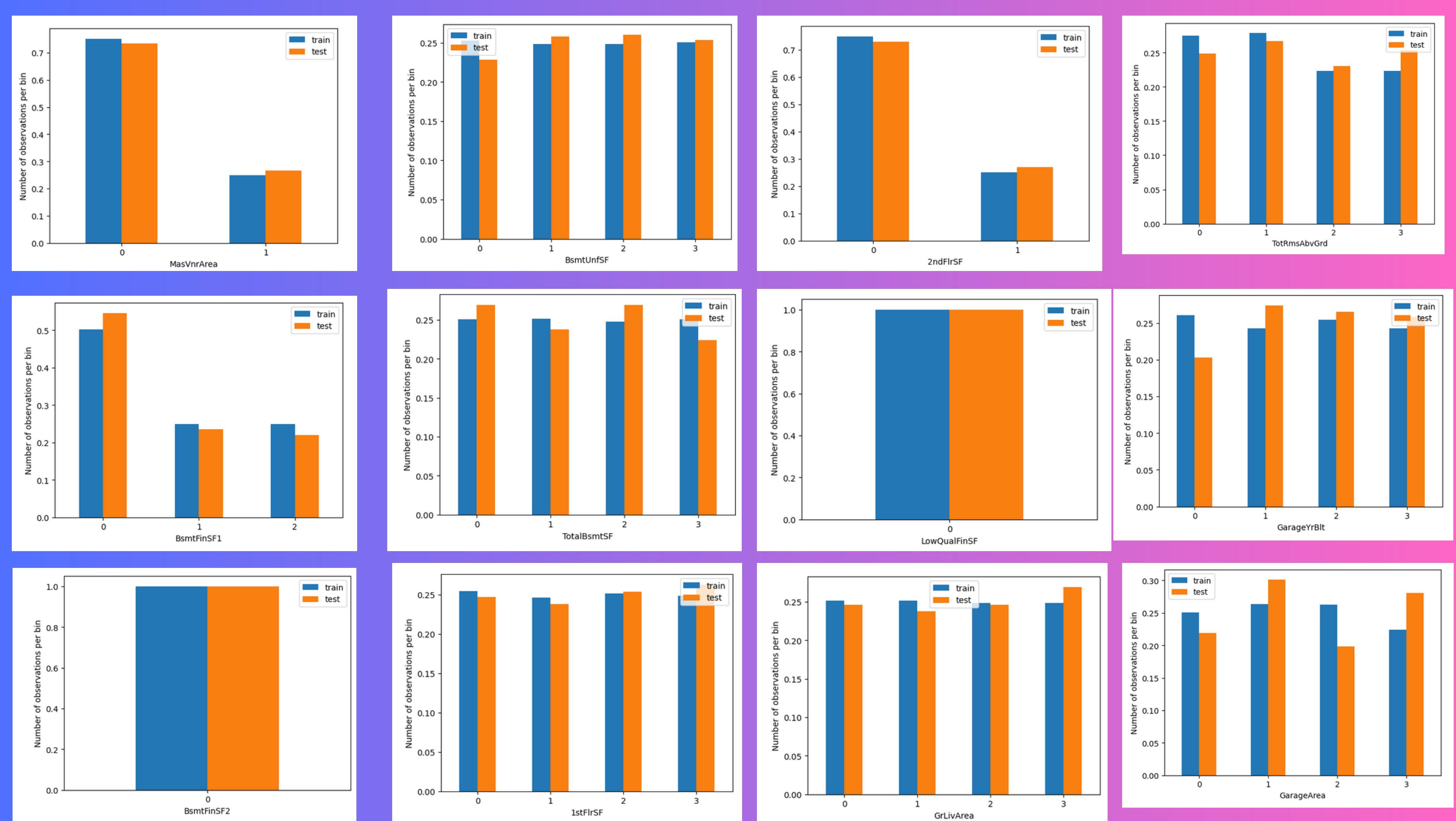
## Equal-frequency

```
1 equal_frequency = EqualFrequencyDiscretiser(q = 4, variables = continuous_features_train)
2
3 equal_frequency.fit(X_train_impute_enc)
4
5 X_train_equal_frequency = equal_frequency.transform(X_train_impute_enc)
6
7 X_test_equal_frequency = equal_frequency.transform(X_test_impute_enc)
```

```
for var in continuous_features_train:
    check_balance_of_discrete_and_categorical(
        X_train_equal_frequency,
        X_test_equal_frequency,
        var
    )
```

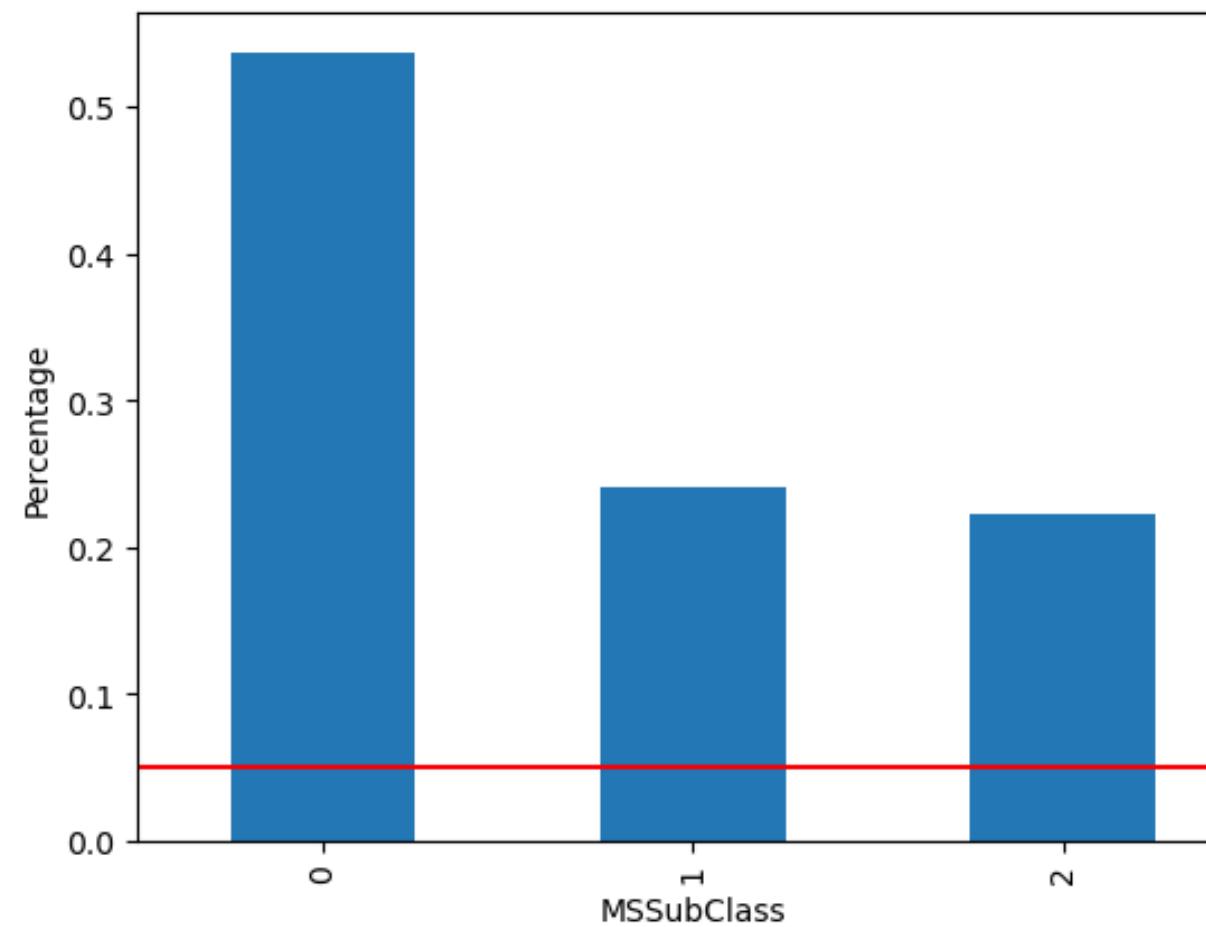
## Check balancing of data between train and test sets





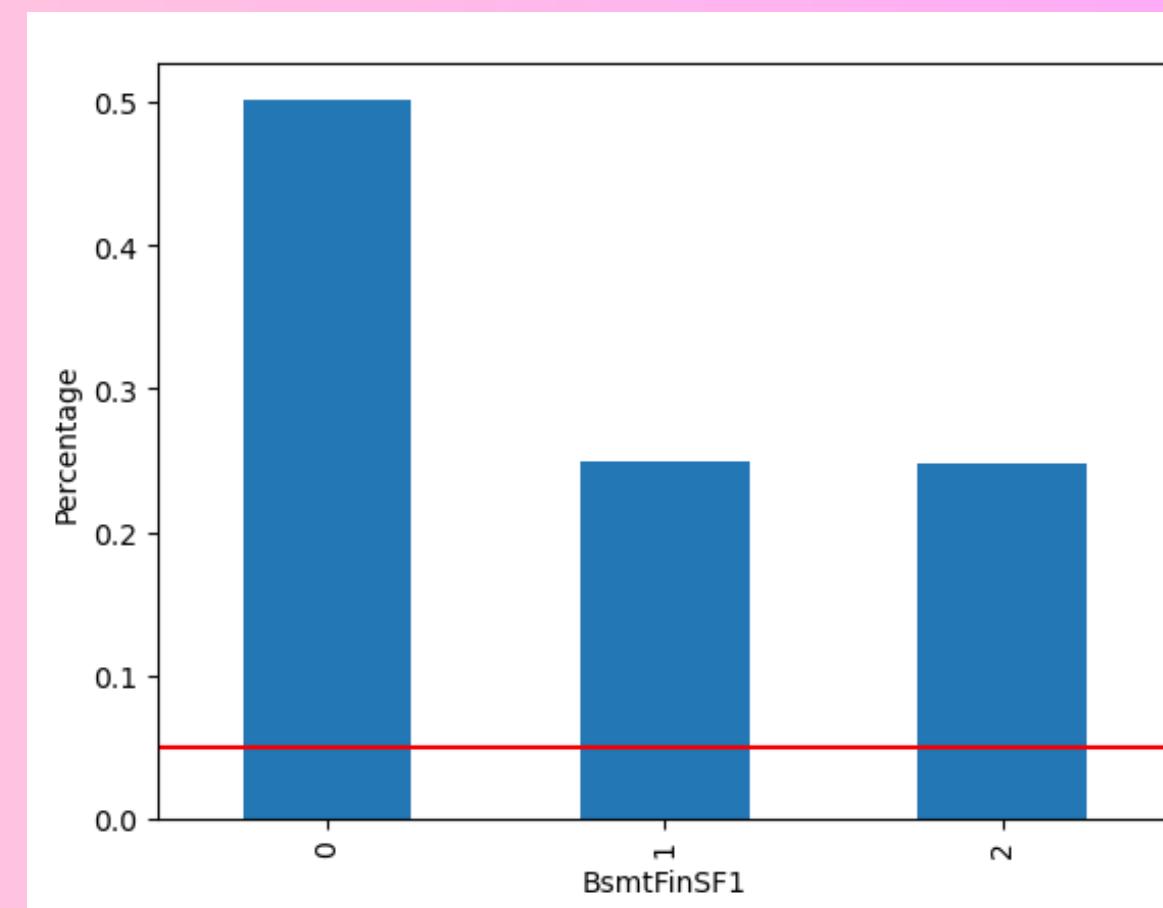
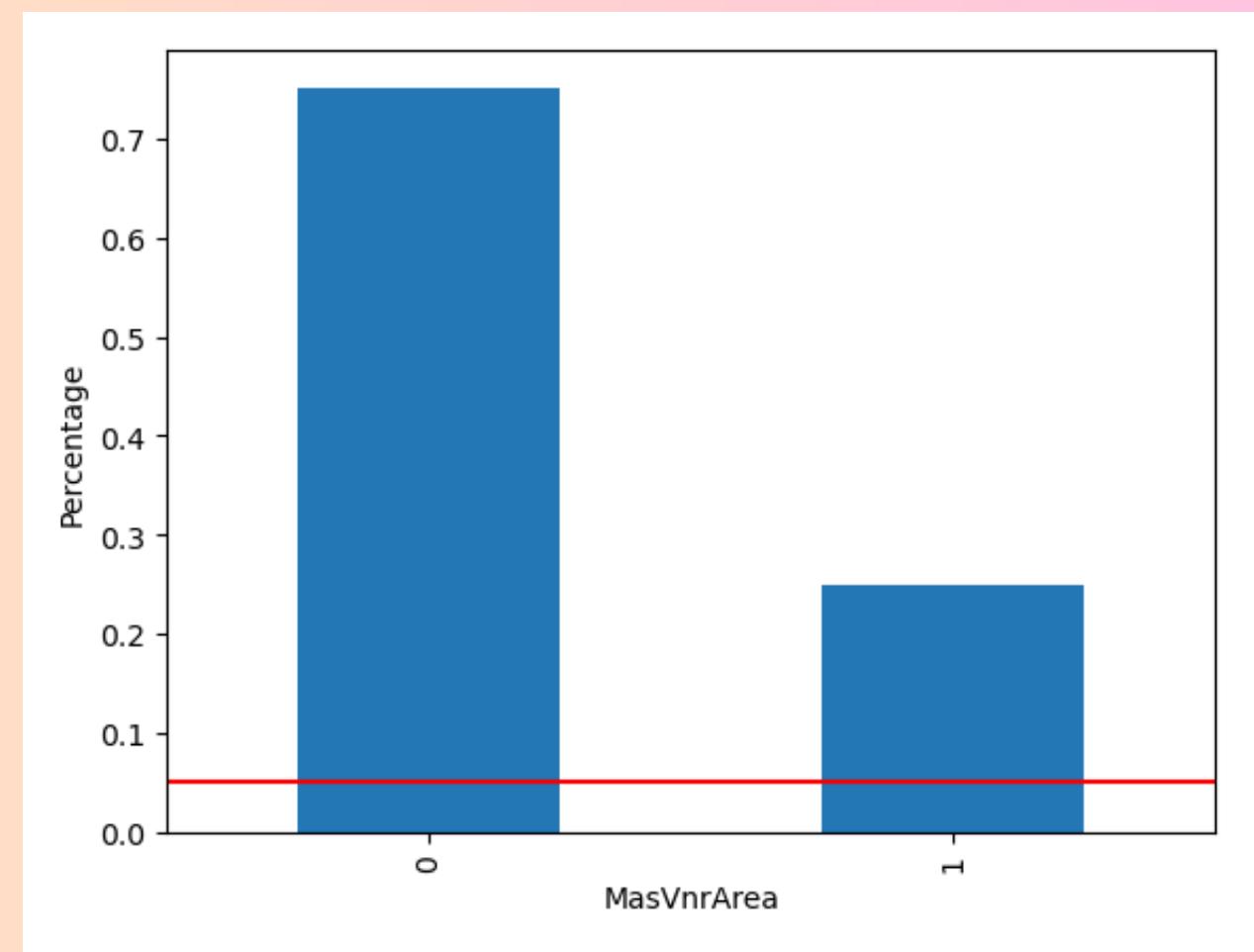
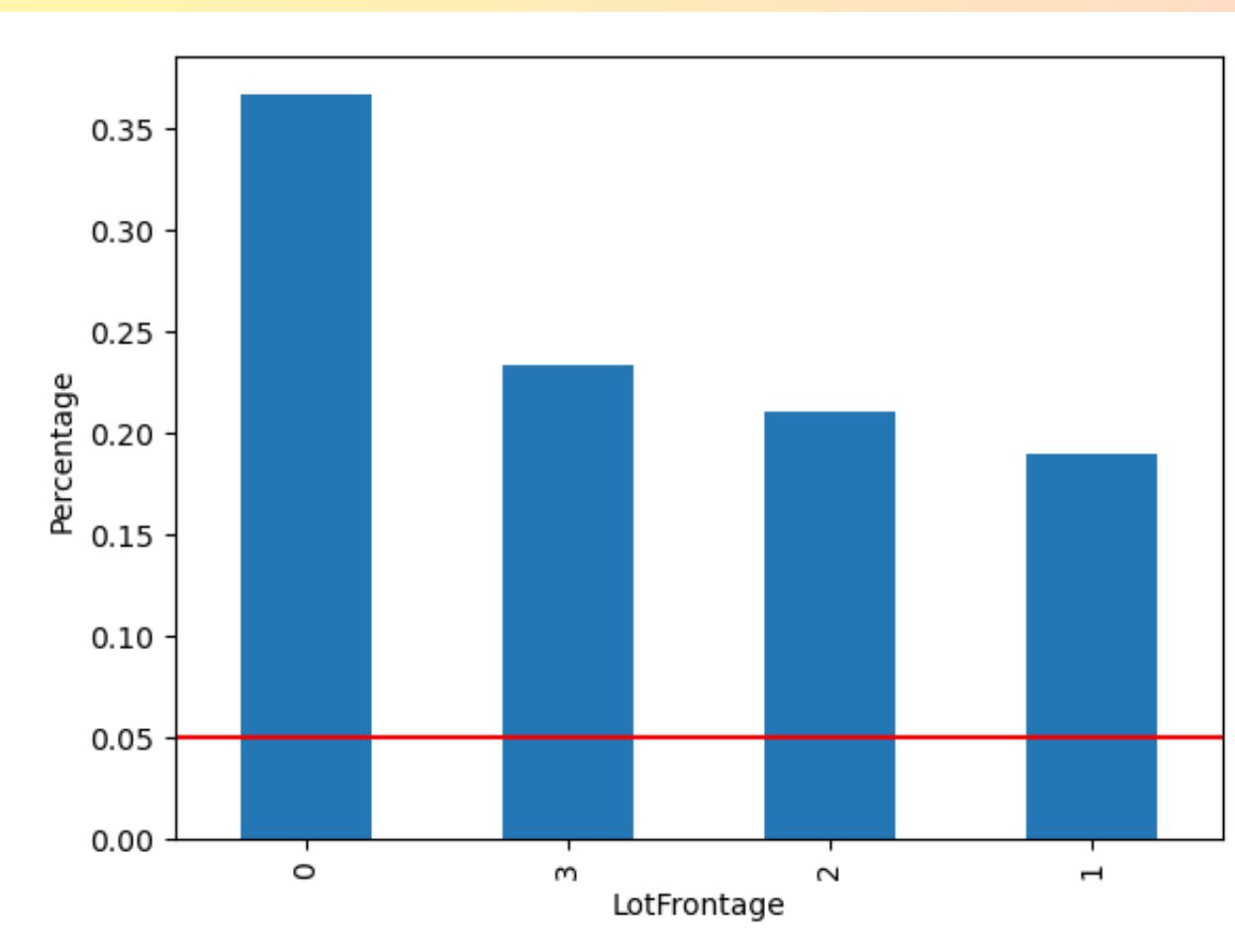
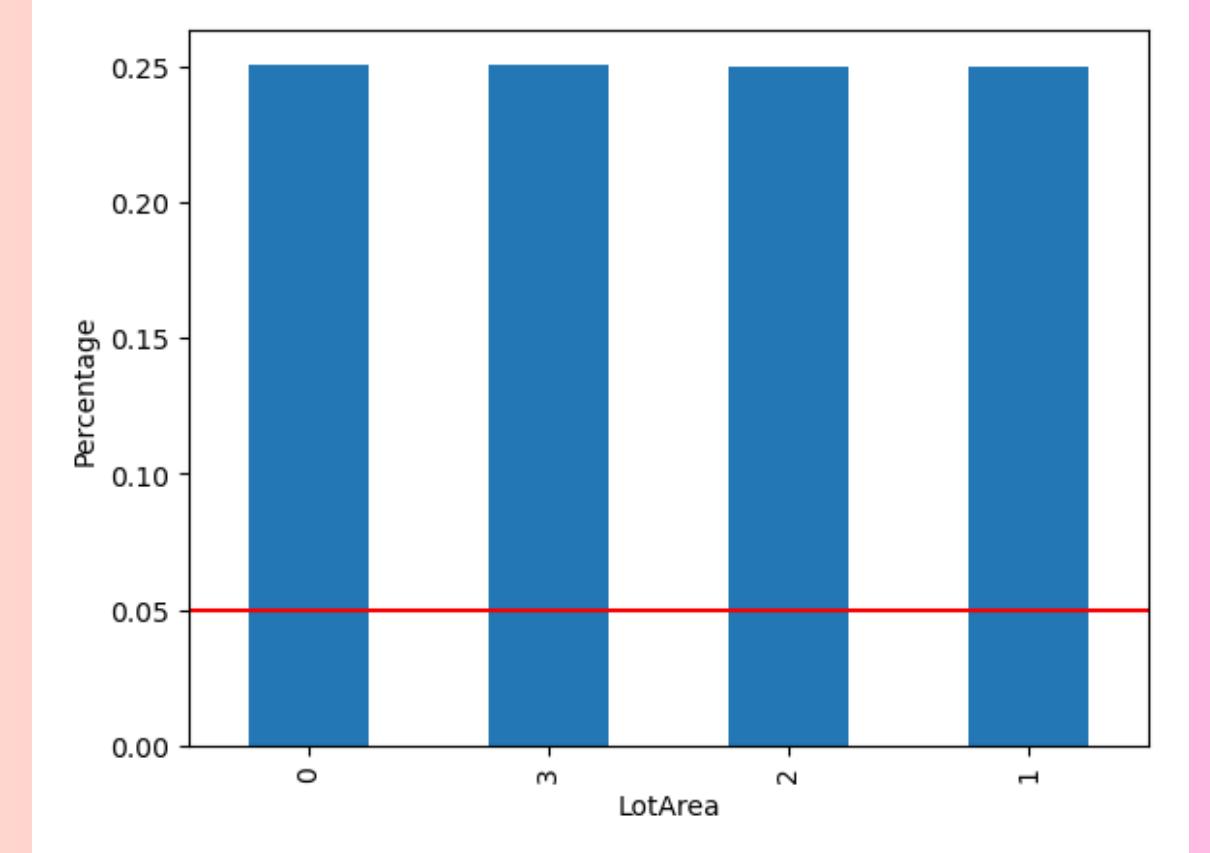
### Rare labels analysis

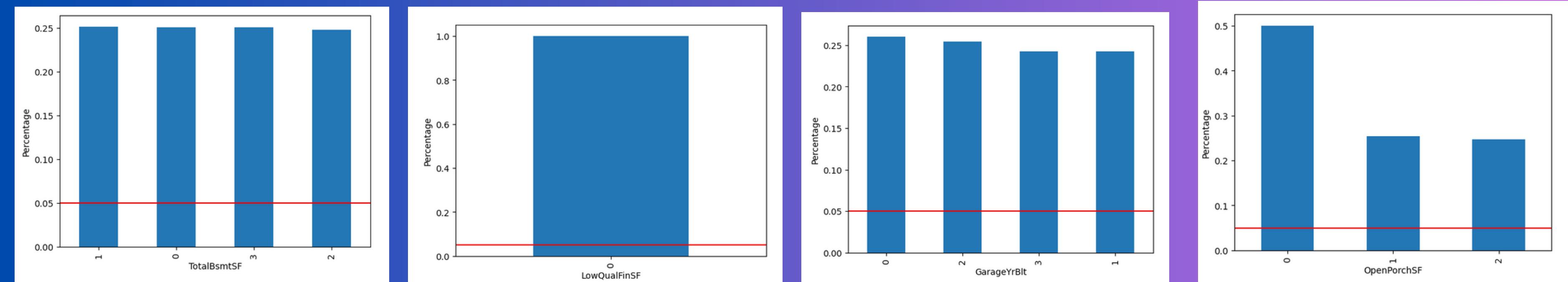
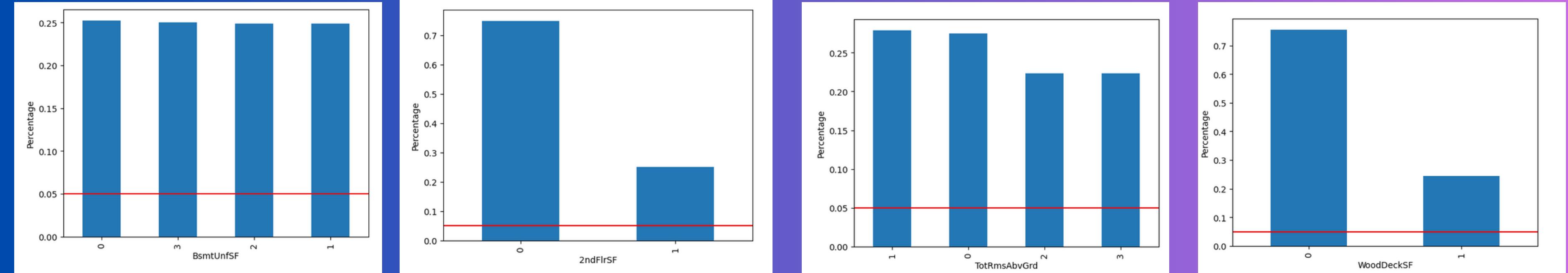
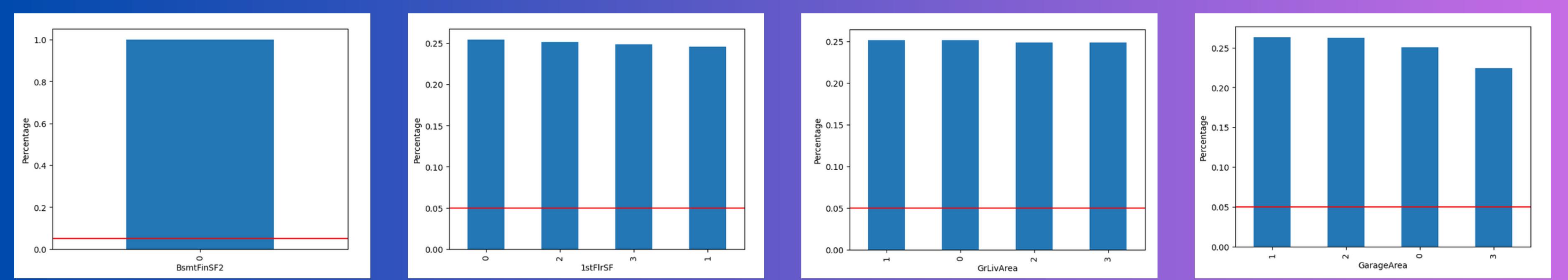
```
: 1 | check_rare_labels_analysis(X_train_equal_frequency, continous_features_train)
```



# Rare label Analysis

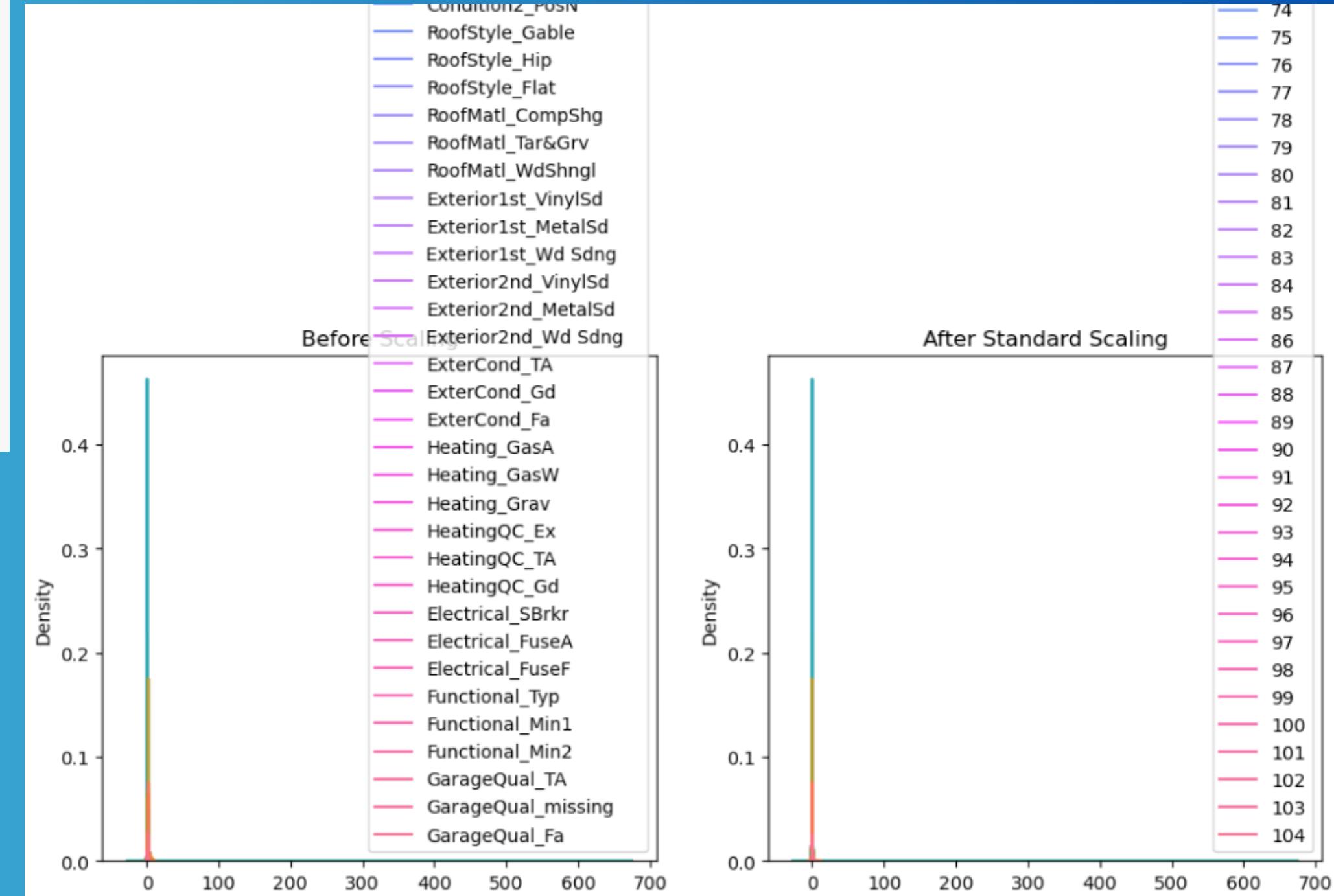
Check rare label for  
X\_train\_equal\_frequency





# Feature Magnitude Robust Scaling

```
scaler_robust = RobustScaler()  
  
scaler_robust.fit(X_train_equal_frequency)  
  
X_train_scaled_robust = scaler_robust.transform(X_train_equal_frequency)  
X_test_scaled_robust = scaler_robust.transform(X_test_equal_frequency)  
  
fig, (ax1, ax2) = plt.subplots(ncols = 2, figsize =(12, 5))  
  
#before scaling  
  
ax1.set_title('Before Scaling')  
sns.kdeplot(X_train_equal_frequency, ax = ax1, warn_singular = False)  
sns.kdeplot(X_train_equal_frequency, ax = ax1, warn_singular = False)  
sns.kdeplot(X_train_equal_frequency, ax = ax1, warn_singular = False)  
  
#after Scaling  
ax2.set_title('After Standard Scaling')  
sns.kdeplot(X_train_scaled_robust, ax = ax2, warn_singular = False)  
sns.kdeplot(X_train_scaled_robust, ax = ax2, warn_singular = False)  
sns.kdeplot(X_train_scaled_robust, ax = ax2, warn_singular = False)  
  
plt.show()
```



# Unit-Norm

```
scaler_unit_norm = Normalizer()
scaler_unit_norm.fit(X_train_equal_frequency)

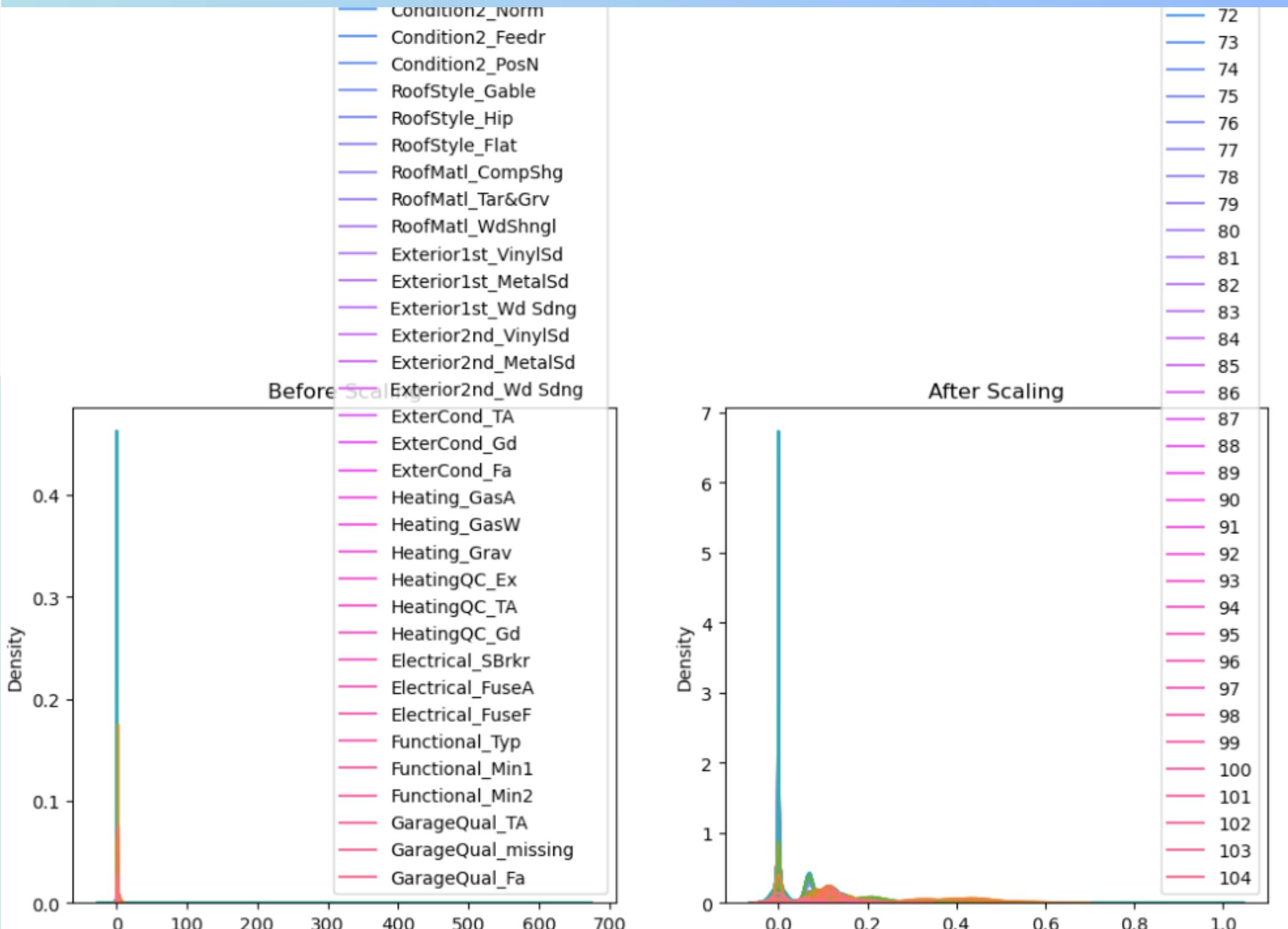
X_train_scaled_unit_norm = scaler_unit_norm.transform(X_train_equal_frequency)

fig, (ax1, ax2) = plt.subplots(ncols = 2, figsize =(12, 5))

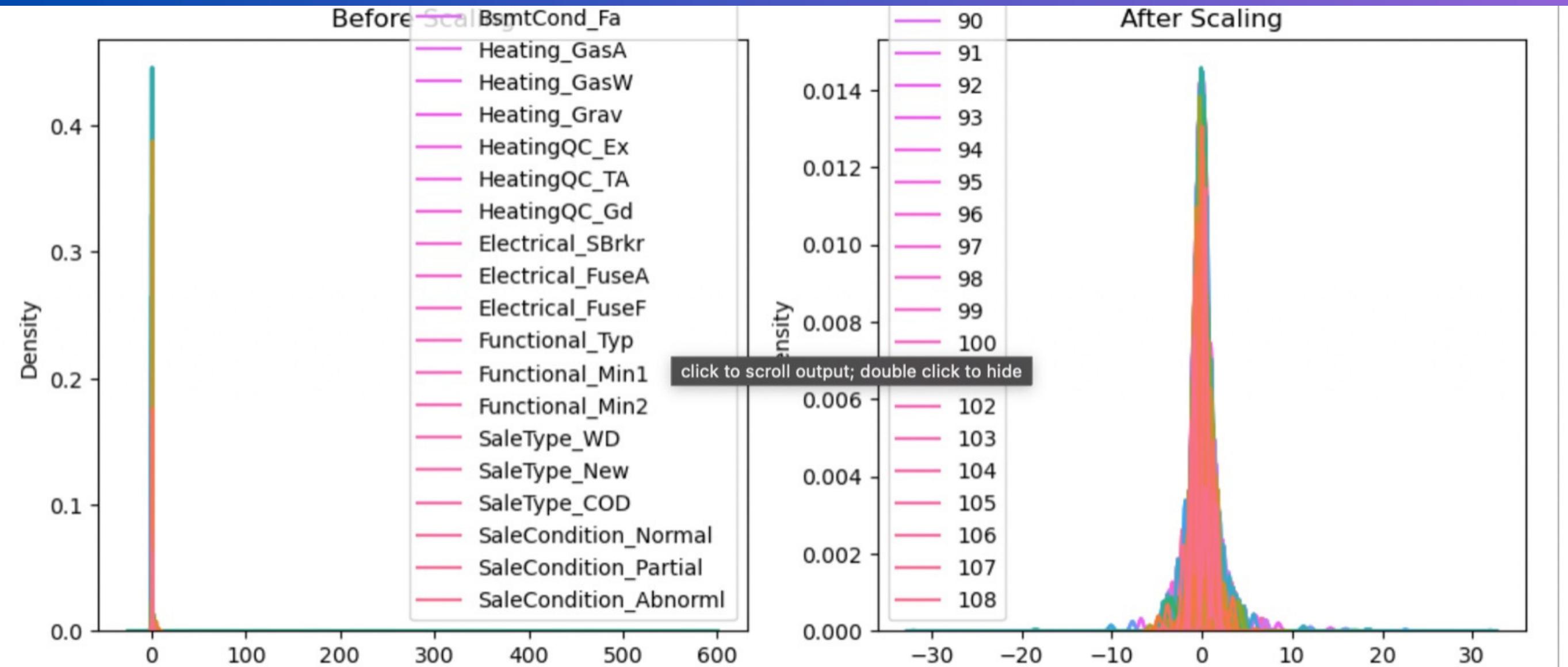
#before scaling
ax1.set_title('Before Scaling')
sns.kdeplot(X_train_equal_frequency, ax = ax1, warn_singular = False)
sns.kdeplot(X_train_equal_frequency, ax = ax1, warn_singular = False)
sns.kdeplot(X_train_equal_frequency, ax = ax1, warn_singular = False)

#after Scaling
ax2.set_title('After Scaling')
sns.kdeplot(X_train_scaled_unit_norm, ax = ax2, warn_singular = False)
sns.kdeplot(X_train_scaled_unit_norm, ax = ax2, warn_singular = False)
sns.kdeplot(X_train_scaled_unit_norm, ax = ax2, warn_singular = False)

plt.show()
```



# Standartization



## 5. Feature-selection

### Constant & Quasi-constant features

```
1 sel_const_variables = DropConstantFeatures(tol = 0.75, variables = None, missing_values = 'raise')
2
3 sel_const_variables.fit(X_train_equal_frequency)
```

#### DropConstantFeatures

```
DropConstantFeatures(tol=0.75)
```

```
1 sel_const_variables.features_to_drop_
```

```
['MSZoning',
 'LandContour',
 'BldgType',
 'MasVnrArea',
 'BsmtCond',
 'BsmtFinType2',
 'BsmtFinSF2',
 'LowQualFinSF',
 'BsmtHalfBath',
 'KitchenAbvGr',
 'GarageCond',
 'PavedDrive',
 'WoodDeckSF',
 'EnclosedPorch',
 '3SsnPorch',
 'ScreenPorch',
 'PoolArea',
 'MiscVal',
 'SaleType',
 'SaleCondition',
 'Street_Pave',
 'Street_Grvl',
 'Utilities_AllPub',
 'Utilities_NoSeWa',
 'LandSlope_Gtl',
 'LandSlope_Sev',
 'LandSlope_Mod',
 'CentralAir_Y',
 'CentralAir_N',
 'Neighborhood_NAmes',
 'Neighborhood_CollgCr',
```

```
'Neighborhood_OldTown',
 'Condition1_Norm',
 'Condition1_Feedr',
 'Condition1_Artery',
 'Condition2_Norm',
 'Condition2_Feedr',
 'Condition2_PosN',
 'RoofStyle_Gable',
 'RoofStyle_Hip',
 'RoofStyle_Flat',
 'RoofMatl_CompShg',
 'RoofMatl_Tar&Grv',
 'RoofMatl_WdShngl',
 'Exterior1st_MetalSd',
 'Exterior1st_Wd Sdng',
 'Exterior2nd_MetalSd',
 'Exterior2nd_Wd Sdng',
 'ExterCond_TA',
 'ExterCond_Gd',
 'ExterCond_Fa',
 'Heating_GasA',
 'Heating_GasW',
 'Heating_Grav',
 'HeatingQC_Gd',
 'Electrical_SBrkr',
 'Electrical_FuseA',
 'Electrical_FuseF',
 'Functional_Typ',
 'Functional_Min1',
 'Functional_Min2',
 'GarageQual_TA',
 'GarageQual_missing',
 'GarageQual_Fa']
```

# Feature Selection Constant and Quasi- Constant Features

## Duplicated features

```
1 sel_duplicated_variables = DropDuplicateFeatures(variables = None, missing_values = 'raise')
2
3 sel_duplicated_variables.fit(X_train_equal_frequency)
```

```
▼     DropDuplicateFeatures
```

```
DropDuplicateFeatures(missing_values='raise')
```

```
1 sel_duplicated_variables.features_to_drop_
```

```
{'3SsnPorch', 'EnclosedPorch', 'LowQualFinSF', 'MiscVal', 'ScreenPorch'}
```

# Duplicated Features

## Correlated features

```
1 sel_correlated_variables_variance = SmartCorrelatedSelection(  
2     variables = None,  
3     method = 'pearson',  
4     threshold = 0.8,  
5     missing_values = 'raise',  
6     selection_method = 'variance',  
7     estimator = None,  
8     scoring = 'accuracy',  
9     cv = 5  
10)  
11  
12 sel_correlated_variables_variance.fit(X_train_equal_frequency)
```

SmartCorrelatedSelection

```
SmartCorrelatedSelection(cv=5, missing_values='raise', scoring='accuracy',  
selection_method='variance')
```

# Correlated Features

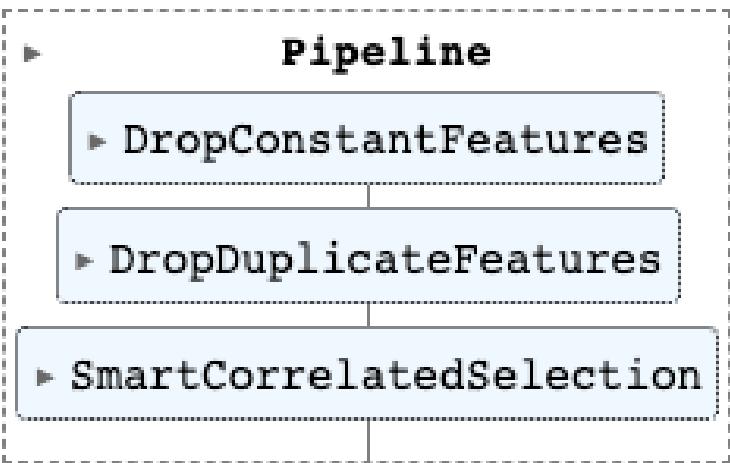
```
1 sel_correlated_variables_variance.features_to_drop_
```

```
[ 'FireplaceQu',  
  'SaleType',  
  'Street_Pave',  
  'Utilities_AllPub',  
  'LandSlope_Mod',  
  'CentralAir_Y',  
  'RoofStyle_Hip',  
  'RoofMatl_Tar&Grv',  
  'Exterior2nd_VinylSd',  
  'Exterior2nd_MetalSd',  
  'Exterior2nd_Wd_Sdng',  
  'ExterCond_Gd',  
  'Electrical_FuseA',  
  'GarageQual_missing' ]
```

```

1 pipe_basic_filter = Pipeline(steps = [
2     ('const_and_quasi-const', DropConstantFeatures(tol = 0.8, variables = None, missing_values = 'raise')),
3     ('duplicated', DropDuplicateFeatures(variables = None, missing_values = 'raise')),
4     ('correlation', SmartCorrelatedSelection(
5         variables = None,
6         method = 'pearson',
7         threshold = 0.8,
8         missing_values = 'raise',
9         selection_method = 'variance',
10        estimator = None,
11        scoring = 'accuracy',
12        cv = 5))
13    ])
14
15 pipe_basic_filter.fit(X_train_equal_frequency)

```



```

1 X_train_basic_filtered = pipe_basic_filter.transform(X_train_equal_frequency)
2 X_test_basic_filtered = pipe_basic_filter.transform(X_test_equal_frequency)

```

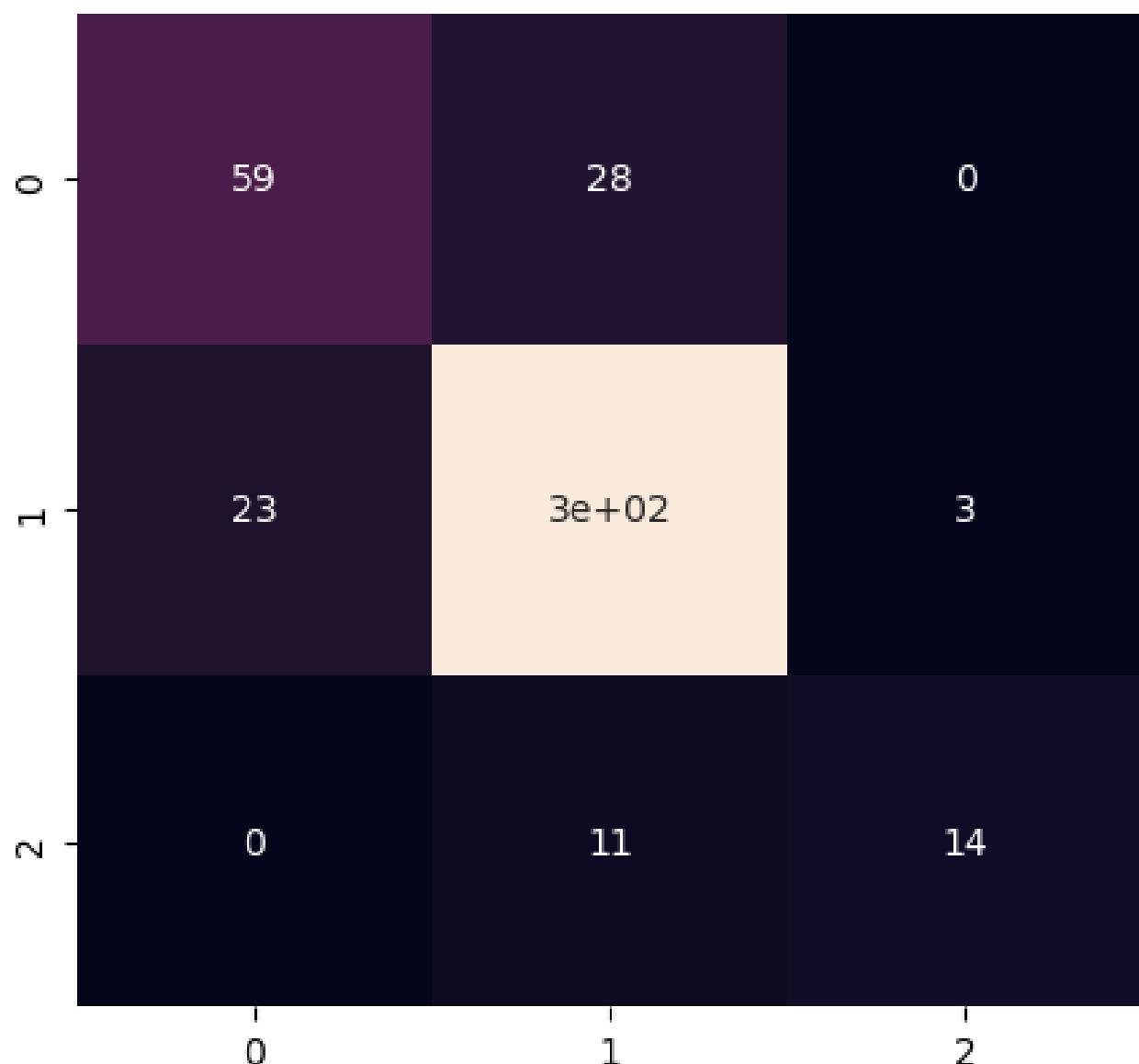
```
1 X_train_basic_filtered.shape, X_test_basic_filtered.shape
```

```
((1022, 43), (438, 43))
```

# Build basic filter methods with PipeLine

# Metrics Score

## Decision Tree



Classification report on train set

	precision	recall	f1-score	support
0	0.72	0.68	0.70	87
1	0.88	0.92	0.90	326
2	0.82	0.56	0.67	25
accuracy			0.85	438
macro avg	0.81	0.72	0.76	438
weighted avg	0.85	0.85	0.85	438



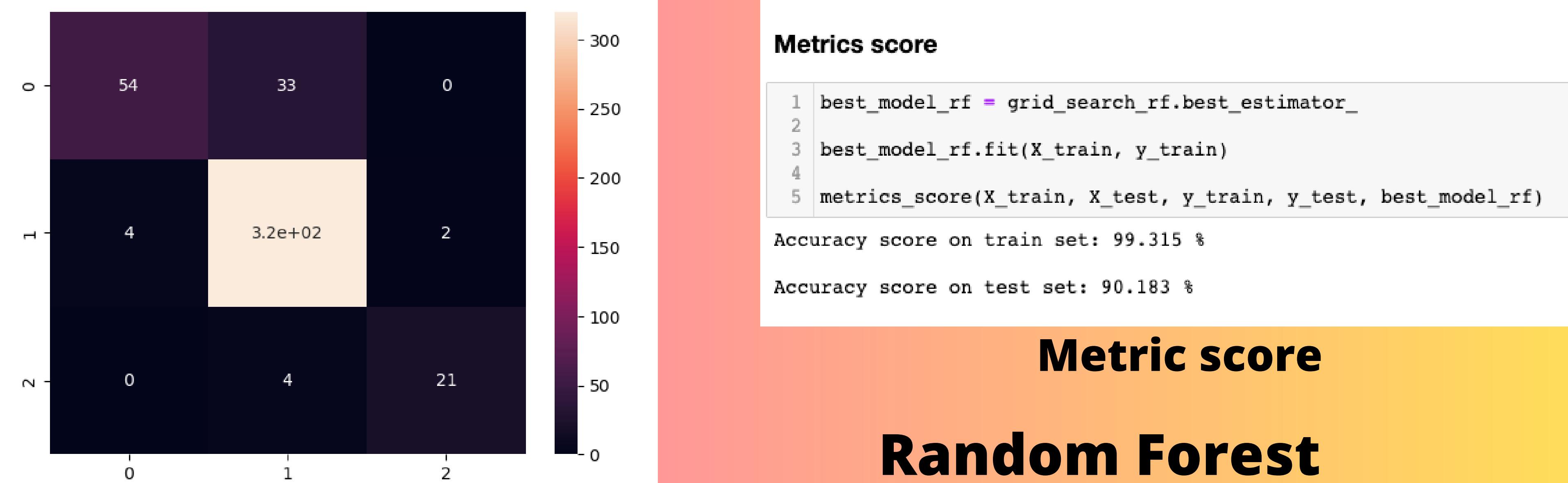
```
: 1 best_model_dt = grid_search_dt.best_estimator_
 2
 3 best_model_dt.fit(X_train, y_train)
 4
 5 metrics_score(X_train, X_test, y_train, y_test, best_model_dt)
```

Accuracy score on train set: 94.814 %

Accuracy score on test set: 85.160 %

```
1 | grid_search_dt.best_params_
```

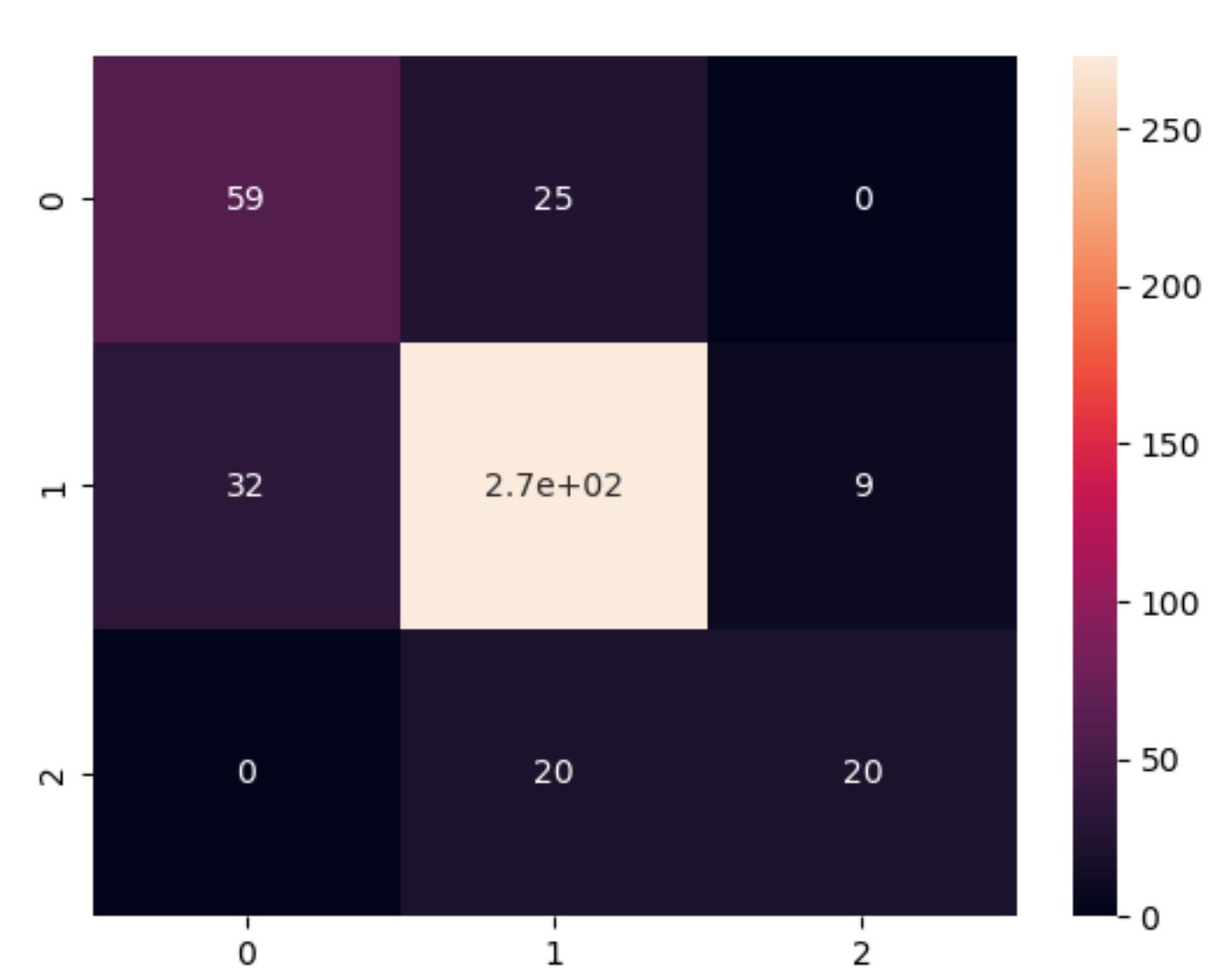
```
('classifier_criterion': 'gini',
 'classifier_max_depth': 10,
 'classifier_max_features': 'log2',
 'classifier_splitter': 'best',
 'disc_q': 6,
 'median_imputation_method': 'median',
 'ohe_top_cat_top_categories': 3,
 'random_sample_random_state': 10}
```



Classification report on train set

	precision	recall	f1-score	support
0	0.93	0.62	0.74	87
1	0.90	0.98	0.94	326
2	0.91	0.84	0.87	25
accuracy			0.90	438
macro avg	0.91	0.81	0.85	438
weighted avg	0.90	0.90	0.90	438

```
1 grid_search_rf.best_params_
{'classifier_criterion': 'log_loss',
 'classifier_max_depth': 9,
 'classifier_n_estimators': 200,
 'disc_q': 8,
 'median_imputation_method': 'median'}
```



Classification report on train set

	precision	recall	f1-score	support
0	0.65	0.70	0.67	84
1	0.86	0.87	0.86	314
2	0.69	0.50	0.58	40
accuracy			0.80	438
macro avg	0.73	0.69	0.71	438
weighted avg	0.80	0.80	0.80	438

```
1 grid_search_ada.best_params_
{'classifier_learning_rate': 0.5,
'classifier_n_estimators': 100,
'disc_q': 8,
'median_imputation_method': 'median',
'ohe_top_cat_top_categories': 3}
```

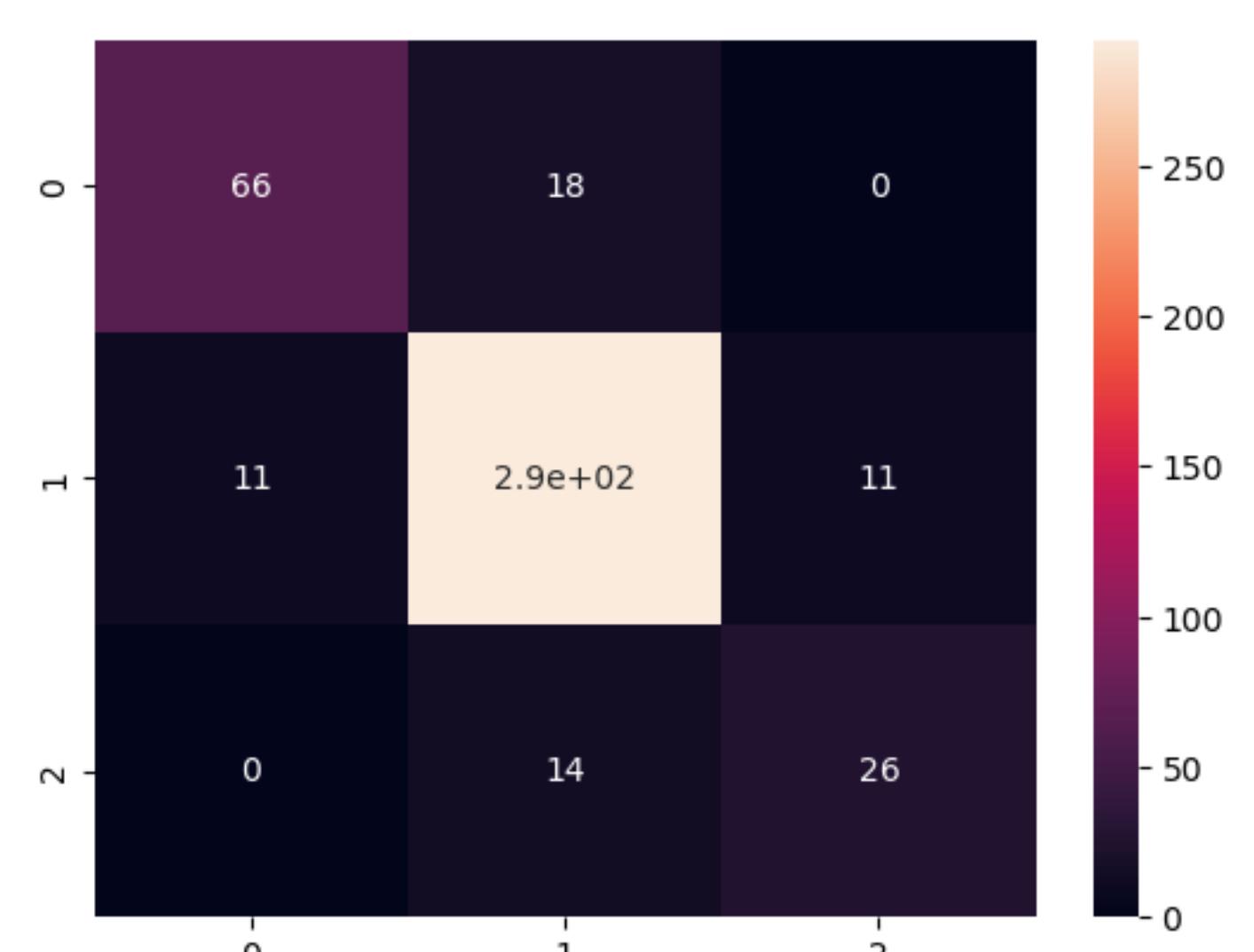
## Metric score

# Ada Boost Classifier

```
1 best_model_ada = grid_search_ada.best_estimator_
2
3 best_model_ada.fit(X_train, y_train)
4
5 metrics_score(X_train, X_test, y_train, y_test, best_model_ada)
```

Accuracy score on train set: 80.920 %

Accuracy score on test set: 80.365 %



Classification report on train set

	precision	recall	f1-score	support
0	0.86	0.79	0.82	84
1	0.90	0.93	0.92	314
2	0.70	0.65	0.68	40
accuracy			0.88	438
macro avg	0.82	0.79	0.80	438
weighted avg	0.87	0.88	0.88	438

```
1 best_model_logit = grid_search_logit.best_estimator_
2
3 best_model_logit.fit(X_train, y_train)
4
5 metrics_score(X_train, X_test, y_train, y_test, best_model_logit)
```

Accuracy score on train set: 91.781 %

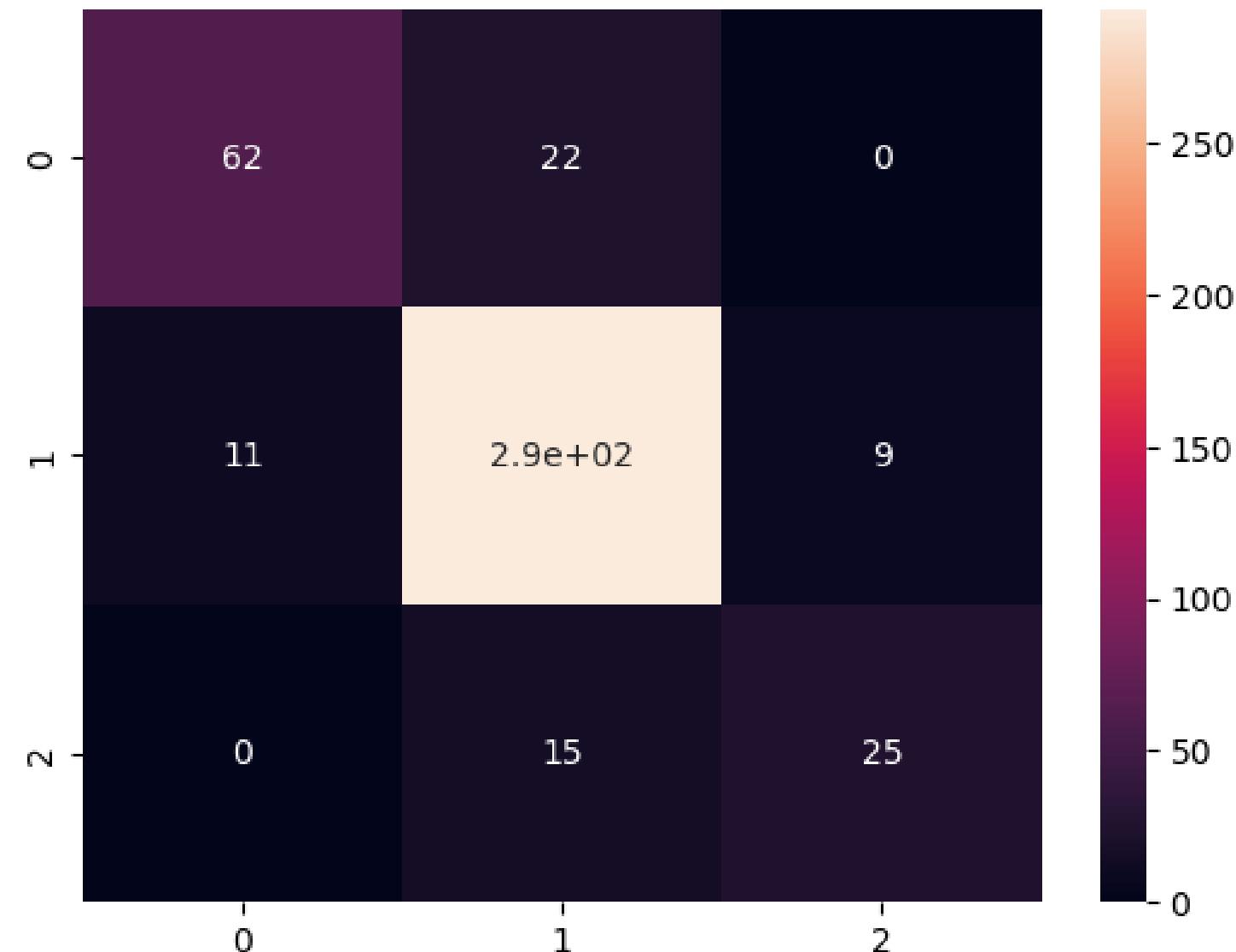
Accuracy score on test set: 87.671 %

## Metric score

## Logistic Regression

```
1 grid_search_logit.best_params_
```

```
{'classifier_C': 0.6,
 'classifier_max_iter': 1000,
 'disc_q': 4,
 'median_imputation_method': 'mean',
 'ohe_top_cat_top_categories': 3}
```



Classification report on train set

	precision	recall	f1-score	support
0	0.85	0.74	0.79	84
1	0.89	0.94	0.91	314
2	0.74	0.62	0.68	40
accuracy			0.87	438
macro avg	0.82	0.77	0.79	438
weighted avg	0.87	0.87	0.87	438

```

1 best_model_svc = grid_search_svc.best_estimator_
2
3 best_model_svc.fit(X_train, y_train)
4
5 metrics_score(X_train, X_test, y_train, y_test, best_model_svc)

```

Accuracy score on train set: 92.074 %

Accuracy score on test set: 86.986 %

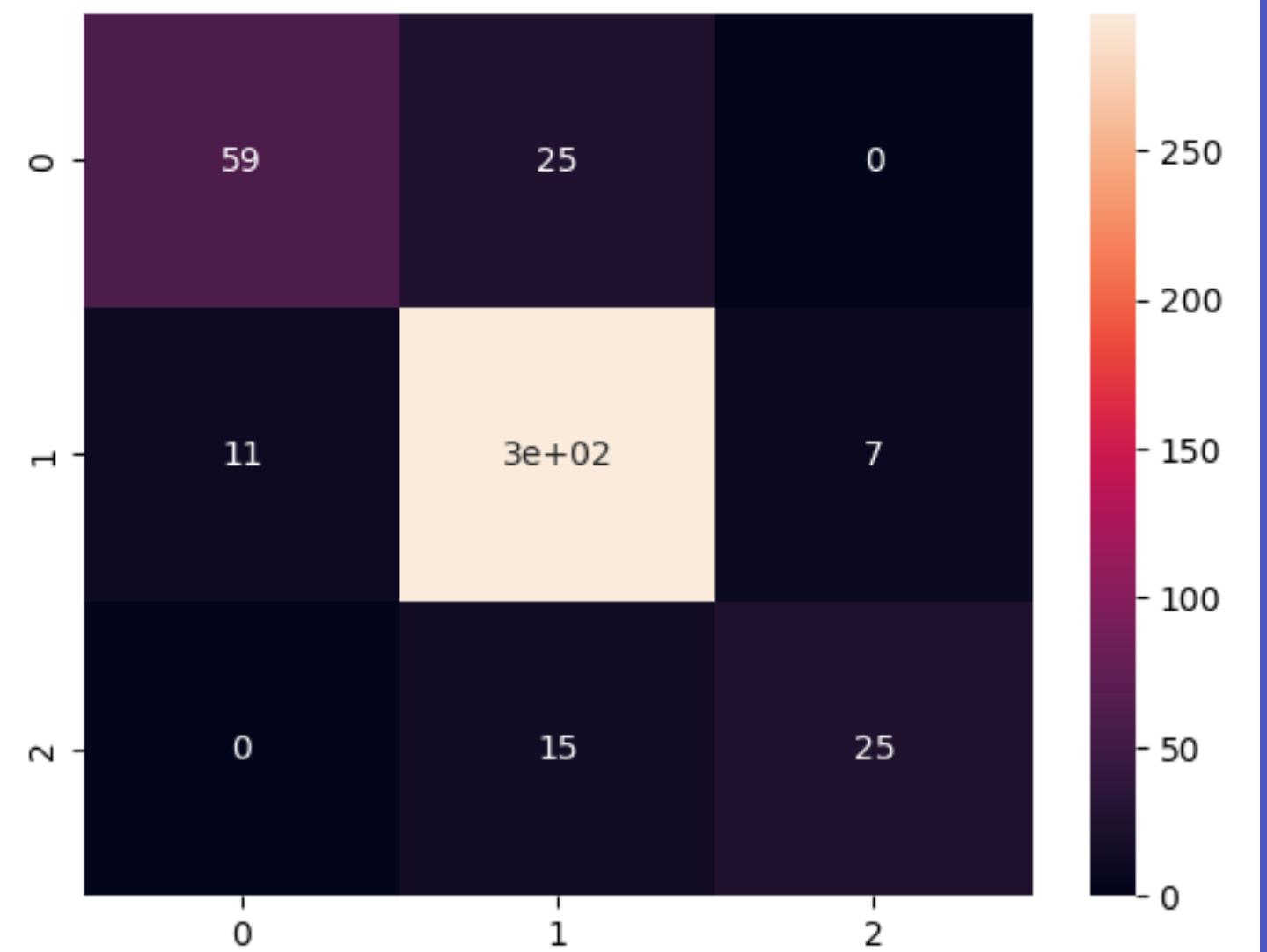
# Metric score

## Support vector classifier

```

1 grid_search_svc.best_params_
{'classifier_C': 0.1,
 'classifier_gamma': 1,
 'classifier_kernel': 'linear',
 'disc_q': 8,
 'median_imputation_method': 'median',
 'ohe_top_cat_top_categories': 3}

```



Classification report on train set

	precision	recall	f1-score	support
0	0.84	0.70	0.77	84
1	0.88	0.94	0.91	314
2	0.78	0.62	0.69	40
accuracy			0.87	438
macro avg	0.84	0.76	0.79	438
weighted avg	0.86	0.87	0.86	438

```

1 grid_search_xgb.best_params_
{'classifier__gamma': 0,
 'classifier__learning_rate': 0.1,
 'classifier__max_depth': 1,
 'classifier__n_estimators': 400,
 'disc_q': 2,
 'median_imputation_method': 'mean',
 'ohe_top_cat_top_categories': 5}

```

# Metric Score

# XGBoost classifier

```

1 best_model_xgb = grid_search_xgb.best_estimator_
2
3 best_model_xgb.fit(X_train, y_train)
4
5 metrics_score(X_train, X_test, y_train, y_test, best_model_xgb)

```

Accuracy score on train set: 91.487 %

Accuracy score on test set: 86.758 %