

Missing values Imputation -Categorical variable

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: df= pd.read_csv("train.csv")
df.head()
```

Out[2]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	...	PoolArea	PoolQC	Fence	Mis
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	

5 rows × 81 columns



```
In [3]: cat_vars=df.select_dtypes(include='object')
cat_vars.head()
```

Out[3]:

	MSZoning	Street	Alley	LotShape	LandContour	Utilities	LotConfig	LandSlope	Neighborhood	Condition1	...	GarageType	Garage
0	RL	Pave	NaN	Reg	Lvl	AllPub	Inside	Gtl	CollgCr	Norm	...	Attchd	
1	RL	Pave	NaN	Reg	Lvl	AllPub	FR2	Gtl	Veenker	Feedr	...	Attchd	
2	RL	Pave	NaN	IR1	Lvl	AllPub	Inside	Gtl	CollgCr	Norm	...	Attchd	
3	RL	Pave	NaN	IR1	Lvl	AllPub	Corner	Gtl	Crawfor	Norm	...	Detchd	
4	RL	Pave	NaN	IR1	Lvl	AllPub	FR2	Gtl	NoRidge	Norm	...	Attchd	

5 rows × 43 columns



```
In [4]: miss_val_per=cat_vars.isnull().mean()*100  
miss_val_per
```

```
Out[4]: MSZoning      0.000000  
Street      0.000000  
Alley      93.767123  
LotShape    0.000000  
LandContour 0.000000  
Utilities   0.000000  
LotConfig   0.000000  
LandSlope   0.000000  
Neighborhood 0.000000  
Condition1  0.000000  
Condition2  0.000000  
BldgType    0.000000  
HouseStyle  0.000000  
RoofStyle   0.000000  
RoofMatl    0.000000  
Exterior1st 0.000000  
Exterior2nd 0.000000  
MasVnrType  0.547945  
ExterQual   0.000000  
ExterCond   0.000000  
Foundation  0.000000  
BsmtQual    2.534247  
BsmtCond    2.534247  
BsmtExposure 2.602740  
BsmtFinType1 2.534247  
BsmtFinType2 2.602740  
Heating      0.000000  
HeatingQC    0.000000  
CentralAir   0.000000  
Electrical   0.068493  
KitchenQual  0.000000  
Functional   0.000000  
FireplaceQu 47.260274  
GarageType   5.547945  
GarageFinish 5.547945
```

```

GarageQual      5.547945
GarageCond      5.547945
PavedDrive      0.000000
PoolQC          99.520548
Fence           80.753425
MiscFeature     96.301370
SaleType        0.000000
SaleCondition   0.000000
dtype: float64

```

```

In [5]: drop_vars=['Alley','FireplaceQu','PoolQC','Fence','MiscFeature']
cat_vars.drop(columns=drop_vars,axis=1,inplace=True)
cat_vars.head()

```

c:\python3.8.3\lib\site-packages\pandas\core\frame.py:4160: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
return super().drop(

Out[5]:

	MSZoning	Street	LotShape	LandContour	Utilities	LotConfig	LandSlope	Neighborhood	Condition1	Condition2	...	Electrical	Kitchen
0	RL	Pave	Reg	Lvl	AllPub	Inside	Gtl	CollgCr	Norm	Norm	...	SBrkr	
1	RL	Pave	Reg	Lvl	AllPub	FR2	Gtl	Veenker	Feedr	Norm	...	SBrkr	
2	RL	Pave	IR1	Lvl	AllPub	Inside	Gtl	CollgCr	Norm	Norm	...	SBrkr	
3	RL	Pave	IR1	Lvl	AllPub	Corner	Gtl	Crawfor	Norm	Norm	...	SBrkr	
4	RL	Pave	IR1	Lvl	AllPub	FR2	Gtl	NoRidge	Norm	Norm	...	SBrkr	

5 rows × 38 columns



```
In [6]: cat_vars.shape
```

```
Out[6]: (1460, 38)
```

```
In [7]: isnull_per = cat_vars.isnull().mean()*100  
miss_vars=isnull_per[isnull_per >0 ].keys()  
miss_vars
```

```
Out[7]: Index(['MasVnrType', 'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1',  
              'BsmtFinType2', 'Electrical', 'GarageType', 'GarageFinish',  
              'GarageQual', 'GarageCond'],  
             dtype='object')
```

```
In [8]: # cat_vars['MasVnrType'].fillna("missing")
```

```
In [9]: cat_vars['MasVnrType'].mode()
```

```
Out[9]: 0      None  
dtype: object
```

```
In [10]: cat_vars['MasVnrType'].value_counts()
```

```
Out[10]: None      864  
BrkFace    445  
Stone      128  
BrkCmn      15  
Name: MasVnrType, dtype: int64
```

```
In [11]: cat_vars['MasVnrType'].fillna(cat_vars['MasVnrType'].mode()[0])
```

```
Out[11]: 0      BrkFace  
1         None  
2      BrkFace  
3         None  
4      BrkFace  
      ...  
1455      None  
1456      Stone  
1457      None  
1458      None  
1459      None  
Name: MasVnrType, Length: 1460, dtype: object
```

```
In [12]: for var in miss_vars:
          cat_vars[var].fillna(cat_vars[var].mode()[0],inplace=True)
          print(var,"=",cat_vars[var].mode()[0])
```

```
MasVnrType = None
BsmtQual = TA
BsmtCond = TA
BsmtExposure = No
BsmtFinType1 = Unf
BsmtFinType2 = Unf
Electrical = SBrkr
GarageType = Attchd
GarageFinish = Unf
GarageQual = TA
GarageCond = TA
```

```
c:\python3.8.3\lib\site-packages\pandas\core\series.py:4517: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
return super().fillna(
```

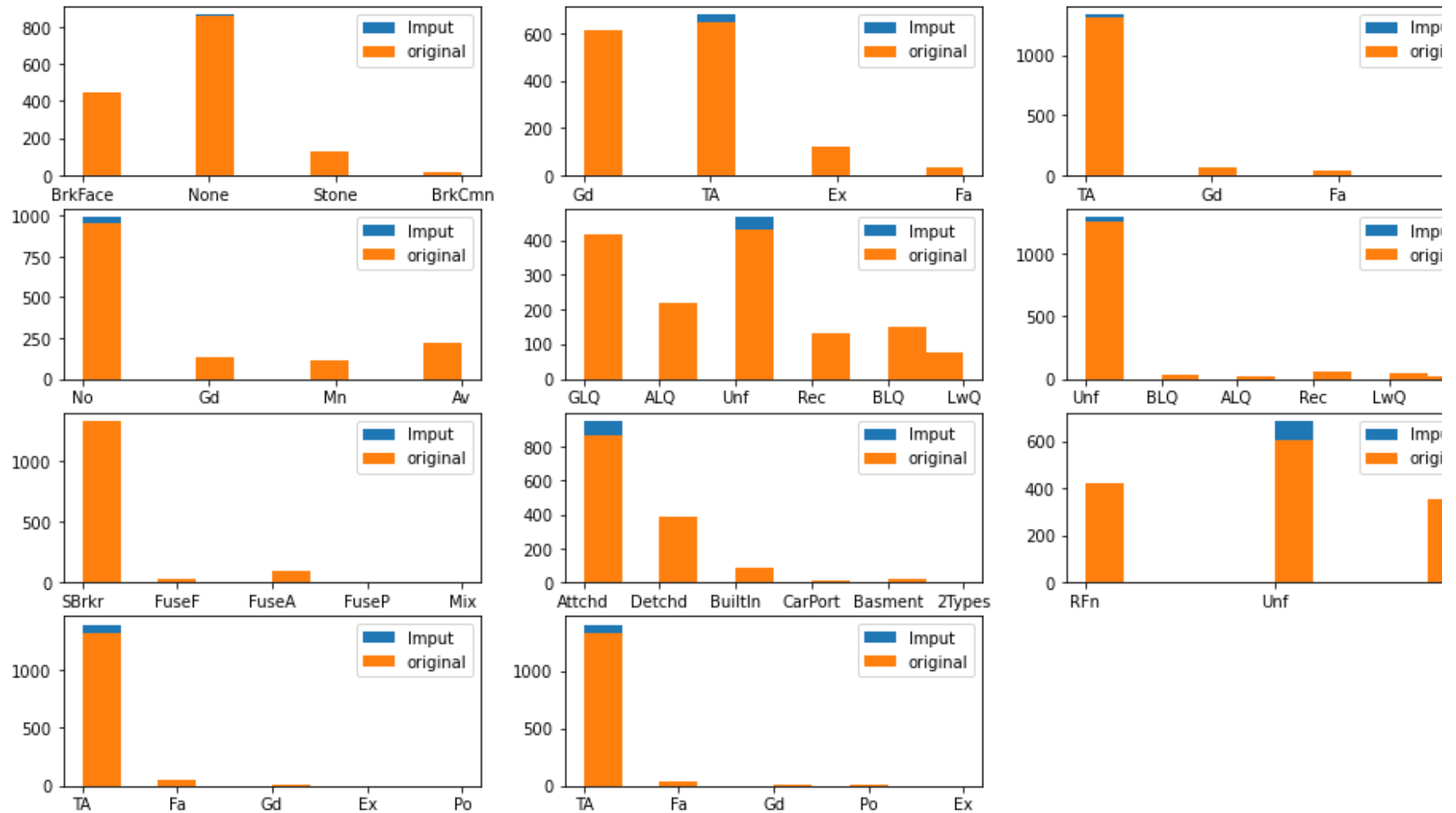
```
In [13]: cat_vars.isnull().sum()
```

```
Out[13]: MSZoning      0  
Street      0  
LotShape    0  
LandContour 0  
Utilities   0  
LotConfig   0  
LandSlope   0  
Neighborhood 0  
Condition1  0  
Condition2  0  
BldgType    0  
HouseStyle  0  
RoofStyle   0  
RoofMat1    0  
Exterior1st 0  
Exterior2nd 0  
MasVnrType  0  
ExterQual   0  
ExterCond   0  
Foundation  0  
BsmtQual    0  
BsmtCond    0  
BsmtExposure 0  
BsmtFinType1 0  
BsmtFinType2 0  
Heating     0  
HeatingQC   0  
CentralAir  0  
Electrical  0  
KitchenQual 0  
Functional  0  
GarageType  0  
GarageFinish 0  
GarageQual  0  
GarageCond  0  
PavedDrive  0
```



```
SaleType      0  
SaleCondition 0  
dtype: int64
```

```
In [14]: plt.figure(figsize=(16,9))
for i,var in enumerate(miss_vars):
    plt.subplot(4,3,i+1)
    plt.hist(cat_vars[var],label= 'Imput')
    plt.hist(df[var].dropna(),label= 'original')
    plt.legend()
```



```
In [15]: df.update(cat_vars)
df.drop(columns=drop_vars,inplace=True)
```

```
In [16]: df.select_dtypes(include='object').isnull().sum()
```

```
Out[16]: MSZoning      0
          Street      0
          LotShape     0
          LandContour  0
          Utilities    0
          LotConfig    0
          LandSlope    0
          Neighborhood  0
          Condition1   0
          Condition2   0
          BldgType      0
          HouseStyle    0
          RoofStyle     0
          RoofMatl      0
          Exterior1st   0
          Exterior2nd   0
          MasVnrType    0
          ExterQual     0
          ExterCond     0
          Foundation    0
          BsmtQual      0
          BsmtCond      0
          BsmtExposure  0
          BsmtFinType1  0
          BsmtFinType2  0
          Heating      0
          HeatingQC     0
          CentralAir    0
          Electrical    0
          KitchenQual   0
          Functional    0
          GarageType    0
          GarageFinish  0
          GarageQual    0
          GarageCond    0
          PavedDrive    0
```

```
SaleType      0  
SaleCondition 0  
dtype: int64
```

