

Missind values imputation using sklearn

for Numeric and Categorical data

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

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

print('train dataset shape :-',train.shape)
print('test dataset shape :-',test.shape)
```

```
train dataset shape :- (1460, 81)
test dataset shape :- (1459, 80)
```

In [3]: `train.head()`

Out[3]:

	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 [4]:

```

x_train=train.drop(columns='SalePrice')
y_train=train['SalePrice']

print('train dataset shape :-',x_train.shape)
print('test dataset shape :-',y_train.shape)

```

train dataset shape :- (1460, 80)
test dataset shape :- (1460,)

numerical values inputation

```
In [13]: num_vars=x_train.select_dtypes(include=["int64","float64"]).columns
num_vars
```

```
Out[13]: Index(['Id', 'MSSubClass', 'LotFrontage', 'LotArea', 'OverallQual',
               'OverallCond', 'YearBuilt', 'YearRemodAdd', 'MasVnrArea', 'BsmtFinSF1',
               'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', '1stFlrSF', '2ndFlrSF',
               'LowQualFinSF', 'GrLivArea', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath',
               'HalfBath', 'BedroomAbvGr', 'KitchenAbvGr', 'TotRmsAbvGrd',
               'Fireplaces', 'GarageYrBlt', 'GarageCars', 'GarageArea', 'WoodDeckSF',
               'OpenPorchSF', 'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'PoolArea',
               'MiscVal', 'MoSold', 'YrSold'],
              dtype='object')
```

```
In [17]: x_train[num_vars].isnull().sum()
```

```
Out[17]: Id                0
MSSubClass                0
LotFrontage              259
LotArea                  0
OverallQual               0
OverallCond               0
YearBuilt                 0
YearRemodAdd              0
MasVnrArea                8
BsmtFinSF1                0
BsmtFinSF2                0
BsmtUnfSF                 0
TotalBsmtSF               0
1stFlrSF                  0
2ndFlrSF                  0
LowQualFinSF              0
GrLivArea                 0
BsmtFullBath              0
BsmtHalfBath              0
FullBath                  0
HalfBath                  0
BedroomAbvGr              0
KitchenAbvGr              0
TotRmsAbvGrd              0
Fireplaces                0
GarageYrBlt              81
GarageCars                0
GarageArea                0
WoodDeckSF                0
OpenPorchSF               0
EnclosedPorch             0
3SsnPorch                 0
ScreenPorch               0
PoolArea                  0
MiscVal                   0
MoSold                    0
```

```
YrSold          0  
dtype: int64
```

```
In [20]: imputer_mean=SimpleImputer(strategy="mean")  
# imputer_mean=SimpleImputer(strategy="constant",fill_value=999....jo valuse tuze impute karni heoo values
```

```
In [22]: imputer_mean.fit(x_train[num_vars])
```

```
Out[22]: SimpleImputer()
```

```
In [23]: imputer_mean.statistics_
```

```
Out[23]: array([7.30500000e+02, 5.68972603e+01, 7.00499584e+01, 1.05168281e+04,  
                6.09931507e+00, 5.57534247e+00, 1.97126781e+03, 1.98486575e+03,  
                1.03685262e+02, 4.43639726e+02, 4.65493151e+01, 5.67240411e+02,  
                1.05742945e+03, 1.16262671e+03, 3.46992466e+02, 5.84452055e+00,  
                1.51546370e+03, 4.25342466e-01, 5.75342466e-02, 1.56506849e+00,  
                3.82876712e-01, 2.86643836e+00, 1.04657534e+00, 6.51780822e+00,  
                6.13013699e-01, 1.97850616e+03, 1.76712329e+00, 4.72980137e+02,  
                9.42445205e+01, 4.66602740e+01, 2.19541096e+01, 3.40958904e+00,  
                1.50609589e+01, 2.75890411e+00, 4.34890411e+01, 6.32191781e+00,  
                2.00781575e+03])
```

```
In [24]: imputer_mean.transform(x_train[num_vars])
```

```
Out[24]: array([[1.000e+00, 6.000e+01, 6.500e+01, ..., 0.000e+00, 2.000e+00,  
                2.008e+03],  
               [2.000e+00, 2.000e+01, 8.000e+01, ..., 0.000e+00, 5.000e+00,  
                2.007e+03],  
               [3.000e+00, 6.000e+01, 6.800e+01, ..., 0.000e+00, 9.000e+00,  
                2.008e+03],  
               ...,  
               [1.458e+03, 7.000e+01, 6.600e+01, ..., 2.500e+03, 5.000e+00,  
                2.010e+03],  
               [1.459e+03, 2.000e+01, 6.800e+01, ..., 0.000e+00, 4.000e+00,  
                2.010e+03],  
               [1.460e+03, 2.000e+01, 7.500e+01, ..., 0.000e+00, 6.000e+00,  
                2.008e+03]])
```

```
In [30]: x_train[num_vars]=imputer_mean.transform(x_train[num_vars])  
         test[num_vars]=imputer_mean.transform(test[num_vars])
```

```
In [32]: x_train[num_vars].isnull().sum()
```

```
Out[32]: Id                0
MSSubClass                0
LotFrontage               0
LotArea                  0
OverallQual               0
OverallCond               0
YearBuilt                 0
YearRemodAdd              0
MasVnrArea                0
BsmtFinSF1                0
BsmtFinSF2                0
BsmtUnfSF                 0
TotalBsmtSF               0
1stFlrSF                  0
2ndFlrSF                  0
LowQualFinSF              0
GrLivArea                 0
BsmtFullBath              0
BsmtHalfBath              0
FullBath                  0
HalfBath                  0
BedroomAbvGr              0
KitchenAbvGr              0
TotRmsAbvGrd              0
Fireplaces                0
GarageYrBlt               0
GarageCars                0
GarageArea                0
WoodDeckSF                0
OpenPorchSF               0
EnclosedPorch             0
3SsnPorch                 0
ScreenPorch               0
PoolArea                  0
MiscVal                   0
MoSold                    0
```

```
YrSold      0  
dtype: int64
```



```
In [34]: test[num_vars].isnull().sum()
```

```
Out[34]: Id                0
MSSubClass                0
LotFrontage               0
LotArea                   0
OverallQual               0
OverallCond               0
YearBuilt                 0
YearRemodAdd              0
MasVnrArea                0
BsmtFinSF1                0
BsmtFinSF2                0
BsmtUnfSF                 0
TotalBsmtSF               0
1stFlrSF                  0
2ndFlrSF                  0
LowQualFinSF              0
GrLivArea                 0
BsmtFullBath              0
BsmtHalfBath              0
FullBath                  0
HalfBath                  0
BedroomAbvGr              0
KitchenAbvGr              0
TotRmsAbvGrd              0
Fireplaces                0
GarageYrBlt               0
GarageCars                0
GarageArea                0
WoodDeckSF                0
OpenPorchSF               0
EnclosedPorch             0
3SsnPorch                 0
ScreenPorch               0
PoolArea                  0
MiscVal                   0
MoSold                    0
```

YrSold 0
dtype: int64

categorical missing values imputation valuse

```
In [37]: cat_vars=x_train.select_dtypes(include=["O"]).columns # objects = 0 ..capiac 0
cat_vars
```

```
Out[37]: Index(['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'],
              dtype='object')
```

```
In [38]: x_train[cat_vars].isnull().sum()
```

```
Out[38]: MSZoning      0
Street      0
Alley      1369
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  8
ExterQual   0
ExterCond   0
Foundation  0
BsmtQual    37
BsmtCond    37
BsmtExposure 38
BsmtFinType1 37
BsmtFinType2 38
Heating     0
HeatingQC   0
CentralAir  0
Electrical  1
KitchenQual 0
Functional  0
FireplaceQu 690
GarageType  81
GarageFinish 81
GarageQual  81
```

```

GarageCond      81
PavedDrive      0
PoolQC          1453
Fence           1179
MiscFeature     1406
SaleType        0
SaleCondition   0
dtype: int64

```

```
In [62]: imputer_mode=SimpleImputer(strategy="most_frequent") #... most_frequent =mode
```

```
In [64]: imputer_mode.fit(x_train[cat_vars])
```

```
Out[64]: SimpleImputer(strategy='most_frequent')
```

```
In [69]: imputer_mode.statistics_
```

```
Out[69]: array(['RL', 'Pave', 'Grvl', 'Reg', 'Lvl', 'AllPub', 'Inside', 'Gtl',
                'NAmes', 'Norm', 'Norm', '1Fam', '1Story', 'Gable', 'CompShg',
                'VinylSd', 'VinylSd', 'None', 'TA', 'TA', 'PConc', 'TA', 'TA',
                'No', 'Unf', 'Unf', 'GasA', 'Ex', 'Y', 'SBrkr', 'TA', 'Typ', 'Gd',
                'Attchd', 'Unf', 'TA', 'TA', 'Y', 'Gd', 'MnPrv', 'Shed', 'WD',
                'Normal'], dtype=object)
```

```
In [70]: imputer_mode.transform(x_train[cat_vars])
```

```
Out[70]: array([[ 'RL', 'Pave', 'Grvl', ..., 'Shed', 'WD', 'Normal'],
                [ 'RL', 'Pave', 'Grvl', ..., 'Shed', 'WD', 'Normal'],
                [ 'RL', 'Pave', 'Grvl', ..., 'Shed', 'WD', 'Normal'],
                ...,
                [ 'RL', 'Pave', 'Grvl', ..., 'Shed', 'WD', 'Normal'],
                [ 'RL', 'Pave', 'Grvl', ..., 'Shed', 'WD', 'Normal'],
                [ 'RL', 'Pave', 'Grvl', ..., 'Shed', 'WD', 'Normal']], dtype=object)
```

```
In [50]: x_train[cat_vars]=imputer_mode.transform(x_train[cat_vars])  
         test[cat_vars]=imputer_mode.transform(test[cat_vars])
```

```
In [51]: x_train[cat_vars].isnull().sum()
```

```
Out[51]: MSZoning      0
Street      0
Alley       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
FireplaceQu  0
GarageType   0
GarageFinish 0
GarageQual   0
```

```
GarageCond      0
PavedDrive      0
PoolQC          0
Fence           0
MiscFeature     0
SaleType        0
SaleCondition   0
dtype: int64
```

```
In [52]: test[cat_vars].isnull().sum()
```

```
Out[52]: MSZoning      0
Street      0
Alley       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
FireplaceQu  0
GarageType   0
GarageFinish 0
GarageQual   0
```



```
GarageCond      0
PavedDrive      0
PoolQC          0
Fence           0
MiscFeature     0
SaleType        0
SaleCondition    0
dtype: int64
```

```
In [53]: x_train[cat_vars].isnull().sum().sum()
```

```
Out[53]: 0
```

```
In [ ]:
```