

# preprocessing

November 9, 2020

```
[1]: import numpy as np
import pandas as pd
```

```
[14]: data = pd.read_csv("data/housing_prices/train.csv")
```

```
[15]: print(data.shape)
data.head()
```

(1460, 81)

```
[15]:
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	\
0	1	60	RL	65.0	8450	Pave	NaN	Reg	
1	2	20	RL	80.0	9600	Pave	NaN	Reg	
2	3	60	RL	68.0	11250	Pave	NaN	IR1	
3	4	70	RL	60.0	9550	Pave	NaN	IR1	
4	5	60	RL	84.0	14260	Pave	NaN	IR1	

  

	LandContour	Utilities	...	PoolArea	PoolQC	Fence	MiscFeature	MiscVal	\
0	Lvl	AllPub	...	0	NaN	NaN	NaN	0	
1	Lvl	AllPub	...	0	NaN	NaN	NaN	0	
2	Lvl	AllPub	...	0	NaN	NaN	NaN	0	
3	Lvl	AllPub	...	0	NaN	NaN	NaN	0	
4	Lvl	AllPub	...	0	NaN	NaN	NaN	0	

  

	MoSold	YrSold	SaleType	SaleCondition	SalePrice
0	2	2008	WD	Normal	208500
1	5	2007	WD	Normal	181500
2	9	2008	WD	Normal	223500
3	2	2006	WD	Abnorml	140000
4	12	2008	WD	Normal	250000

[5 rows x 81 columns]

```
[16]: data.columns
```

```
[16]: Index(['Id', 'MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street',
        'Alley', 'LotShape', 'LandContour', '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')

```

```

[22]: # Filter for null values
vars_with_na = [var for var in data.columns if data[var].isnull().sum() > 0]
print(len(vars_with_na))
print(vars_with_na)

```

19

```

['LotFrontage', 'Alley', 'MasVnrType', 'MasVnrArea', 'BsmtQual', 'BsmtCond',
'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2', 'Electrical', 'FireplaceQu',
'GarageType', 'GarageYrBlt', 'GarageFinish', 'GarageQual', 'GarageCond',
'PoolQC', 'Fence', 'MiscFeature']

```

```

[24]: # print(data[vars_with_na])
# determine percentage of missing values
data[vars_with_na].isnull().mean()

```

```

[24]: LotFrontage    0.177397
      Alley         0.937671
      MasVnrType    0.005479
      MasVnrArea    0.005479
      BsmtQual      0.025342
      BsmtCond      0.025342
      BsmtExposure  0.026027
      BsmtFinType1  0.025342
      BsmtFinType2  0.026027
      Electrical    0.000685
      FireplaceQu   0.472603
      GarageType    0.055479
      GarageYrBlt   0.055479
      GarageFinish  0.055479
      GarageQual    0.055479
      GarageCond    0.055479

```

```
PoolQC      0.995205
Fence       0.807534
MiscFeature  0.963014
dtype: float64
```

```
[ ]:
```