

## Exercise

For this exercise, you will be working with the [House Price Dataset](#).

Please grab the train.csv file from Kaggle and explore this dataset. You need to perform exploratory data analysis and see if there is any correlation between the variables and analyze the distribution of the dataset. The question is open-ended and basically you're asked to perform EDA.

1- Write a summary of your findings in one page (e.g., summary statistics, plots) and submit the pdf file. Therefore, for part 3 of your assignment, you need to submit at least one jupyter notebook file and one pdf file.

2- Push your code and project to github and provide the link to your code here. Ensure that your github project is organized to at least couple of main folders, ensure that you have the README file as well:

- Src
- Data
- Docs
- Results

Read this link for further info:

<https://gist.github.com/ericmlj/27e50331f24db3e8f957d1fe7bbbe510>

## Overview

There are **1460** instances of training data

## Read Data

```
In [1]: import pandas as pd  
df_ss=pd.read_csv('..\Data\sample_submission.csv')  
df_test=pd.read_csv('..\Data\test.csv')  
df_train=pd.read_csv('..\Data\train.csv')
```

```
In [2]: df_ss.head()
```

Out[2]:

|          | <b>Id</b> | <b>SalePrice</b> |
|----------|-----------|------------------|
| <b>0</b> | 1461      | 169277.052498    |
| <b>1</b> | 1462      | 187758.393989    |
| <b>2</b> | 1463      | 183583.683570    |
| <b>3</b> | 1464      | 179317.477511    |
| <b>4</b> | 1465      | 150730.079977    |

In [3]:

```
df_test.head()
```

Out[3]:

|          | <b>Id</b> | <b>MSSubClass</b> | <b>MSZoning</b> | <b>LotFrontage</b> | <b>LotArea</b> | <b>Street</b> | <b>Alley</b> | <b>LotShape</b> | <b>LandContour</b> | <b>Utilit</b> |
|----------|-----------|-------------------|-----------------|--------------------|----------------|---------------|--------------|-----------------|--------------------|---------------|
| <b>0</b> | 1461      | 20                | RH              | 80.0               | 11622          | Pave          | NaN          | Reg             |                    | Lvl AllP      |
| <b>1</b> | 1462      | 20                | RL              | 81.0               | 14267          | Pave          | NaN          | IR1             |                    | Lvl AllP      |
| <b>2</b> | 1463      | 60                | RL              | 74.0               | 13830          | Pave          | NaN          | IR1             |                    | Lvl AllP      |
| <b>3</b> | 1464      | 60                | RL              | 78.0               | 9978           | Pave          | NaN          | IR1             |                    | Lvl AllP      |
| <b>4</b> | 1465      | 120               | RL              | 43.0               | 5005           | Pave          | NaN          | IR1             |                    | HLS AllP      |

5 rows × 80 columns

In [4]:

```
df_train.head()
```

Out[4]:

|          | <b>Id</b> | <b>MSSubClass</b> | <b>MSZoning</b> | <b>LotFrontage</b> | <b>LotArea</b> | <b>Street</b> | <b>Alley</b> | <b>LotShape</b> | <b>LandContour</b> | <b>Utilities</b> |
|----------|-----------|-------------------|-----------------|--------------------|----------------|---------------|--------------|-----------------|--------------------|------------------|
| <b>0</b> | 1         | 60                | RL              | 65.0               | 8450           | Pave          | NaN          | Reg             |                    | Lvl AllPub       |
| <b>1</b> | 2         | 20                | RL              | 80.0               | 9600           | Pave          | NaN          | Reg             |                    | Lvl AllPub       |
| <b>2</b> | 3         | 60                | RL              | 68.0               | 11250          | Pave          | NaN          | IR1             |                    | Lvl AllPub       |
| <b>3</b> | 4         | 70                | RL              | 60.0               | 9550           | Pave          | NaN          | IR1             |                    | Lvl AllPub       |
| <b>4</b> | 5         | 60                | RL              | 84.0               | 14260          | Pave          | NaN          | IR1             |                    | Lvl AllPub       |

5 rows × 81 columns

To get access to all columns including the dependent variable for prediction **SalePrice**, the dataframe chose to this EDA is *df\_train*

## General review

First, it is important to know the data, including its **data type** of all dataframe attributes and the number of **missing values** present

## Data Type

*Before checking the missing values, it is important identify if the data type are well defined regarding its content*

In [5]: `df_train.info()`

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1460 entries, 0 to 1459
Data columns (total 81 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   Id                1460 non-null    int64  
 1   MSSubClass         1460 non-null    int64  
 2   MSZoning          1460 non-null    object  
 3   LotFrontage        1201 non-null    float64 
 4   LotArea            1460 non-null    int64  
 5   Street             1460 non-null    object  
 6   Alley              91  non-null     object  
 7   LotShape            1460 non-null    object  
 8   LandContour        1460 non-null    object  
 9   Utilities           1460 non-null    object  
 10  LotConfig           1460 non-null    object  
 11  LandSlope           1460 non-null    object  
 12  Neighborhood        1460 non-null    object  
 13  Condition1          1460 non-null    object  
 14  Condition2          1460 non-null    object  
 15  BldgType            1460 non-null    object  
 16  HouseStyle          1460 non-null    object  
 17  OverallQual         1460 non-null    int64  
 18  OverallCond         1460 non-null    int64  
 19  YearBuilt            1460 non-null    int64  
 20  YearRemodAdd        1460 non-null    int64  
 21  RoofStyle            1460 non-null    object  
 22  RoofMatl             1460 non-null    object  
 23  Exterior1st          1460 non-null    object  
 24  Exterior2nd          1460 non-null    object  
 25  MasVnrType           1452 non-null    object  
 26  MasVnrArea           1452 non-null    float64 
 27  ExterQual            1460 non-null    object  
 28  ExterCond            1460 non-null    object  
 29  Foundation           1460 non-null    object  
 30  BsmtQual             1423 non-null    object  
 31  BsmtCond             1423 non-null    object  
 32  BsmtExposure         1422 non-null    object  
 33  BsmtFinType1          1423 non-null    object  
 34  BsmtFinSF1            1460 non-null    int64  
 35  BsmtFinType2          1422 non-null    object  
 36  BsmtFinSF2            1460 non-null    int64  
 37  BsmtUnfSF             1460 non-null    int64  
 38  TotalBsmtSF           1460 non-null    int64  
 39  Heating               1460 non-null    object  
 40  HeatingQC              1460 non-null    object  
 41  CentralAir            1460 non-null    object  
 42  Electrical            1459 non-null    object  
 43  1stFlrSF              1460 non-null    int64  
 44  2ndFlrSF              1460 non-null    int64  
 45  LowQualFinSF           1460 non-null    int64  
 46  GrLivArea              1460 non-null    int64  
 47  BsmtFullBath           1460 non-null    int64  
 48  BsmtHalfBath           1460 non-null    int64  
 49  FullBath               1460 non-null    int64  
 50  HalfBath                1460 non-null    int64  
 51  BedroomAbvGr            1460 non-null    int64  
 52  KitchenAbvGr            1460 non-null    int64  
 53  KitchenQual             1460 non-null    object  
 54  TotRmsAbvGrd            1460 non-null    int64

```

```

55 Functional      1460 non-null   object
56 Fireplaces       1460 non-null   int64
57 FireplaceQu      770 non-null   object
58 GarageType        1379 non-null   object
59 GarageYrBlt      1379 non-null   float64
60 GarageFinish      1379 non-null   object
61 GarageCars        1460 non-null   int64
62 GarageArea        1460 non-null   int64
63 GarageQual        1379 non-null   object
64 GarageCond        1379 non-null   object
65 PavedDrive        1460 non-null   object
66 WoodDeckSF        1460 non-null   int64
67 OpenPorchSF       1460 non-null   int64
68 EnclosedPorch     1460 non-null   int64
69 3SsnPorch         1460 non-null   int64
70 ScreenPorch        1460 non-null   int64
71 PoolArea          1460 non-null   int64
72 PoolQC            7 non-null     object
73 Fence              281 non-null   object
74 MiscFeature        54 non-null    object
75 MiscVal            1460 non-null   int64
76 MoSold             1460 non-null   int64
77 YrSold             1460 non-null   int64
78 SaleType           1460 non-null   object
79 SaleCondition      1460 non-null   object
80 SalePrice          1460 non-null   int64
dtypes: float64(3), int64(35), object(43)
memory usage: 924.0+ KB

```

*There are inconsistence in the columns types for example:*

- All columns that refers to Year are int type. Even though they are numbers, it does not make sense if sum 2013+2016 or 2013\*2016, we won't do calculation with this numbers. In order how to identify all over the columns if they are set in a aproppiated data type is check one by one.

**There are inconsistencies in the data type of the dataframe columns:**

*For example:*

- All "Year" columns are Integer type.
- **OverallCond** is integer type in the dataset since its content are numbers, but it is actually categories.

*Even though the column content are numbers, it does not mean they are.*

**For example:** There is no sense if we do these calculations:

2013 + 2016 or 2013 \* 2016

We do not do these calculation with Years. So, in order to identify all over the columns with its appropriated data type, the unique values needs to be identify, as it's shown in the following code:

```
In [6]: d={}
for i in df_train.columns:
    d.update({i:(len(df_train[i].unique()), list(df_train[i].unique()))})
```

```
In [7]: pd.set_option('display.max_rows', None, 'display.max_columns', None)
df_unique=pd.DataFrame(d).transpose().sort_values(by=0, ascending=False)
#df_unique[df_unique.Len<5]
df_unique.columns=['Len','Categories']
df_unique
```

Out[7]:

|                      | <b>Len</b> | <b>Categories</b>                                  |
|----------------------|------------|--|
| <b>Id</b>            | 1460       | [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14...] |
| <b>LotArea</b>       | 1073       | [8450, 9600, 11250, 9550, 14260, 14115, 10084,...] |
| <b>GrLivArea</b>     | 861        | [1710, 1262, 1786, 1717, 2198, 1362, 1694, 209...] |
| <b>BsmtUnfSF</b>     | 780        | [150, 284, 434, 540, 490, 64, 317, 216, 952, 1...  |
| <b>1stFlrSF</b>      | 753        | [856, 1262, 920, 961, 1145, 796, 1694, 1107, 1...  |
| <b>TotalBsmtSF</b>   | 721        | [856, 1262, 920, 756, 1145, 796, 1686, 1107, 9...  |
| <b>SalePrice</b>     | 663        | [208500, 181500, 223500, 140000, 250000, 14300...  |
| <b>BsmtFinSF1</b>    | 637        | [706, 978, 486, 216, 655, 732, 1369, 859, 0, 8...  |
| <b>GarageArea</b>    | 441        | [548, 460, 608, 642, 836, 480, 636, 484, 468, ...] |
| <b>2ndFlrSF</b>      | 417        | [854, 0, 866, 756, 1053, 566, 983, 752, 1142, ...] |
| <b>MasVnrArea</b>    | 328        | [196.0, 0.0, 162.0, 350.0, 186.0, 240.0, 286.0...  |
| <b>WoodDeckSF</b>    | 274        | [0, 298, 192, 40, 255, 235, 90, 147, 140, 160,...] |
| <b>OpenPorchSF</b>   | 202        | [61, 0, 42, 35, 84, 30, 57, 204, 4, 21, 33, 21...  |
| <b>BsmtFinSF2</b>    | 144        | [0, 32, 668, 486, 93, 491, 506, 712, 362, 41, ...] |
| <b>EnclosedPorch</b> | 120        | [0, 272, 228, 205, 176, 87, 172, 102, 37, 144,...] |
| <b>YearBuilt</b>     | 112        | [2003, 1976, 2001, 1915, 2000, 1993, 2004, 197...  |
| <b>LotFrontage</b>   | 111        | [65.0, 80.0, 68.0, 60.0, 84.0, 85.0, 75.0, nan...] |
| <b>GarageYrBlt</b>   | 98         | [2003.0, 1976.0, 2001.0, 1998.0, 2000.0, 1993....] |
| <b>ScreenPorch</b>   | 76         | [0, 176, 198, 291, 252, 99, 184, 168, 130, 142...  |
| <b>YearRemodAdd</b>  | 61         | [2003, 1976, 2002, 1970, 2000, 1995, 2005, 197...  |
| <b>Neighborhood</b>  | 25         | [CollgCr, Veenker, Crawfor, NoRidge, Mitchel, ...] |
| <b>LowQualFinSF</b>  | 24         | [0, 360, 513, 234, 528, 572, 144, 392, 371, 39...  |
| <b>MiscVal</b>       | 21         | [0, 700, 350, 500, 400, 480, 450, 15500, 1200,...] |
| <b>3SsnPorch</b>     | 20         | [0, 320, 407, 130, 180, 168, 140, 508, 238, 24...  |
| <b>Exterior2nd</b>   | 16         | [VinylSd, MetalSd, Wd Shng, HdBoard, Plywood, ...] |
| <b>MSSubClass</b>    | 15         | [60, 20, 70, 50, 190, 45, 90, 120, 30, 85, 80,...] |
| <b>Exterior1st</b>   | 15         | [VinylSd, MetalSd, Wd Sdng, HdBoard, BrkFace, ...] |
| <b>MoSold</b>        | 12         | [2, 5, 9, 12, 10, 8, 11, 4, 1, 7, 3, 6]            |
| <b>TotRmsAbvGrd</b>  | 12         | [8, 6, 7, 9, 5, 11, 4, 10, 12, 3, 2, 14]           |
| <b>OverallQual</b>   | 10         | [7, 6, 8, 5, 9, 4, 10, 3, 1, 2]                    |
| <b>SaleType</b>      | 9          | [WD, New, COD, ConLD, ConLI, CWD, ConLw, Con, ...] |
| <b>OverallCond</b>   | 9          | [5, 8, 6, 7, 4, 2, 3, 9, 1]                        |
| <b>Condition1</b>    | 9          | [Norm, Feedr, PosN, Artery, RRAe, RRNn, RRAn, ...] |

|                      | <i>Len</i> | <i>Categories</i>                                   |
|----------------------|------------|---|
| <b>BedroomAbvGr</b>  | 8          | [3, 4, 1, 2, 0, 5, 6, 8]                            |
| <b>PoolArea</b>      | 8          | [0, 512, 648, 576, 555, 480, 519, 738]              |
| <b>RoofMatl</b>      | 8          | [CompShg, WdShngl, Metal, WdShake, Membran, Ta...]  |
| <b>Condition2</b>    | 8          | [Norm, Artery, RRNn, Feedr, PosN, PosA, RRAn, ...]  |
| <b>HouseStyle</b>    | 8          | [2Story, 1Story, 1.5Fin, 1.5Unf, SFoyer, SLvl,...]  |
| <b>BsmtFinType2</b>  | 7          | [Unf, BLQ, nan, ALQ, Rec, LwQ, GLQ]                 |
| <b>GarageType</b>    | 7          | [Attchd, Detachd, BuiltIn, CarPort, nan, Basmen...] |
| <b>BsmtFinType1</b>  | 7          | [GLQ, ALQ, Unf, Rec, BLQ, nan, LwQ]                 |
| <b>Functional</b>    | 7          | [Typ, Min1, Maj1, Min2, Mod, Maj2, Sev]             |
| <b>Electrical</b>    | 6          | [SBrkr, FuseF, FuseA, FuseP, Mix, nan]              |
| <b>RoofStyle</b>     | 6          | [Gable, Hip, Gambrel, Mansard, Flat, Shed]          |
| <b>Heating</b>       | 6          | [GasA, GasW, Grav, Wall, OthW, Floor]               |
| <b>GarageQual</b>    | 6          | [TA, Fa, Gd, nan, Ex, Po]                           |
| <b>Foundation</b>    | 6          | [PConc, CBlock, BrkTil, Wood, Slab, Stone]          |
| <b>SaleCondition</b> | 6          | [Normal, Abnorml, Partial, AdjLand, Alloca, Fa...]  |
| <b>FireplaceQu</b>   | 6          | [nan, TA, Gd, Fa, Ex, Po]                           |
| <b>GarageCond</b>    | 6          | [TA, Fa, nan, Gd, Po, Ex]                           |
| <b>MiscFeature</b>   | 5          | [nan, Shed, Gar2, Othr, TenC]                       |
| <b>Fence</b>         | 5          | [nan, MnPrv, GdWo, GdPrv, MnWw]                     |
| <b>GarageCars</b>    | 5          | [2, 3, 1, 0, 4]                                     |
| <b>YrSold</b>        | 5          | [2008, 2007, 2006, 2009, 2010]                      |
| <b>HeatingQC</b>     | 5          | [Ex, Gd, TA, Fa, Po]                                |
| <b>LotConfig</b>     | 5          | [Inside, FR2, Corner, CulDSac, FR3]                 |
| <b>BldgType</b>      | 5          | [1Fam, 2fmCon, Duplex, TwnhsE, Twnhs]               |
| <b>MasVnrType</b>    | 5          | [BrkFace, None, Stone, BrkCmn, nan]                 |
| <b>ExterCond</b>     | 5          | [TA, Gd, Fa, Po, Ex]                                |
| <b>MSZoning</b>      | 5          | [RL, RM, C (all), FV, RH]                           |
| <b>BsmtQual</b>      | 5          | [Gd, TA, Ex, nan, Fa]                               |
| <b>BsmtCond</b>      | 5          | [TA, Gd, nan, Fa, Po]                               |
| <b>BsmtExposure</b>  | 5          | [No, Gd, Mn, Av, nan]                               |
| <b>KitchenQual</b>   | 4          | [Gd, TA, Ex, Fa]                                    |
| <b>Fireplaces</b>    | 4          | [0, 1, 2, 3]  |
| <b>ExterQual</b>     | 4          | [Gd, TA, Ex, Fa]                                    |
| <b>PoolQC</b>        | 4          | [nan, Ex, Fa, Gd]                                   |

|                     | <b>Len</b> | <b>Categories</b>    |
|---------------------|------------|----------------------|
| <b>LotShape</b>     | 4          | [Reg, IR1, IR2, IR3] |
| <b>LandContour</b>  | 4          | [Lvl, Bnk, Low, HLS] |
| <b>BsmtFullBath</b> | 4          | [1, 0, 2, 3]         |
| <b>GarageFinish</b> | 4          | [RFn, Unf, Fin, nan] |
| <b>FullBath</b>     | 4          | [2, 1, 3, 0]         |
| <b>KitchenAbvGr</b> | 4          | [1, 2, 3, 0]         |
| <b>PavedDrive</b>   | 3          | [Y, N, P]            |
| <b>BsmtHalfBath</b> | 3          | [0, 1, 2]            |
| <b>HalfBath</b>     | 3          | [1, 0, 2]            |
| <b>LandSlope</b>    | 3          | [Gtl, Mod, Sev]      |
| <b>Alley</b>        | 3          | [nan, Grvl, Pave]    |
| <b>Street</b>       | 2          | [Pave, Grvl]         |
| <b>CentralAir</b>   | 2          | [Y, N]               |
| <b>Utilities</b>    | 2          | [AllPub, NoSeWa]     |

*Prior, I have checked the output of the previous table to identify the correct data type. The list that contains this results is saved into datatype variable:*

Let's merge this list into the dataframe `df_unique`

```
In [9]: df_unique['Correct_datatype']=datatype
```

In [10]: `df_unique`

| <i>Out[10]:</i>             | <i>Len</i> |  | <i>Categories</i> | <i>Correct_datatype</i> |
|-----------------------------|------------|--|-------------------|-------------------------|
| <b><i>Id</i></b>            | 1460       | [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14...] |                   | <i>int</i>              |
| <b><i>LotArea</i></b>       | 1073       | [8450, 9600, 11250, 9550, 14260, 14115, 10084...]  |                   | <i>int</i>              |
| <b><i>GrLivArea</i></b>     | 861        | [1710, 1262, 1786, 1717, 2198, 1362, 1694, 209...] |                   | <i>int</i>              |
| <b><i>BsmtUnfSF</i></b>     | 780        | [150, 284, 434, 540, 490, 64, 317, 216, 952, 1...  |                   | <i>int</i>              |
| <b><i>1stFlrSF</i></b>      | 753        | [856, 1262, 920, 961, 1145, 796, 1694, 1107, 1...  |                   | <i>int</i>              |
| <b><i>TotalBsmtSF</i></b>   | 721        | [856, 1262, 920, 756, 1145, 796, 1686, 1107, 9...  |                   | <i>int</i>              |
| <b><i>SalePrice</i></b>     | 663        | [208500, 181500, 223500, 140000, 250000, 14300...  |                   | <i>int</i>              |
| <b><i>BsmtFinSF1</i></b>    | 637        | [706, 978, 486, 216, 655, 732, 1369, 859, 0, 8...  |                   | <i>int</i>              |
| <b><i>GarageArea</i></b>    | 441        | [548, 460, 608, 642, 836, 480, 636, 484, 468, ...] |                   | <i>int</i>              |
| <b><i>2ndFlrSF</i></b>      | 417        | [854, 0, 866, 756, 1053, 566, 983, 752, 1142, ...] |                   | <i>int</i>              |
| <b><i>MasVnrArea</i></b>    | 328        | [196.0, 0.0, 162.0, 350.0, 186.0, 240.0, 286.0...] |                   | <i>float</i>            |
| <b><i>WoodDeckSF</i></b>    | 274        | [0, 298, 192, 40, 255, 235, 90, 147, 140, 160,...] |                   | <i>int</i>              |
| <b><i>OpenPorchSF</i></b>   | 202        | [61, 0, 42, 35, 84, 30, 57, 204, 4, 21, 33, 21...  |                   | <i>int</i>              |
| <b><i>BsmtFinSF2</i></b>    | 144        | [0, 32, 668, 486, 93, 491, 506, 712, 362, 41, ...] |                   | <i>int</i>              |
| <b><i>EnclosedPorch</i></b> | 120        | [0, 272, 228, 205, 176, 87, 172, 102, 37, 144,...] |                   | <i>int</i>              |
| <b><i>YearBuilt</i></b>     | 112        | [2003, 1976, 2001, 1915, 2000, 1993, 2004, 197...  |                   | <i>cat</i>              |
| <b><i>LotFrontage</i></b>   | 111        | [65.0, 80.0, 68.0, 60.0, 84.0, 85.0, 75.0, nan...] |                   | <i>float</i>            |
| <b><i>GarageYrBlt</i></b>   | 98         | [2003.0, 1976.0, 2001.0, 1998.0, 2000.0, 1993....] |                   | <i>cat</i>              |
| <b><i>ScreenPorch</i></b>   | 76         | [0, 176, 198, 291, 252, 99, 184, 168, 130, 142...  |                   | <i>int</i>              |
| <b><i>YearRemodAdd</i></b>  | 61         | [2003, 1976, 2002, 1970, 2000, 1995, 2005, 197...  |                   | <i>cat</i>              |
| <b><i>Neighborhood</i></b>  | 25         | [CollgCr, Veenker, Crawfor, NoRidge, Mitchel, ...] |                   | <i>cat</i>              |
| <b><i>LowQualFinSF</i></b>  | 24         | [0, 360, 513, 234, 528, 572, 144, 392, 371, 39...  |                   | <i>int</i>              |
| <b><i>MiscVal</i></b>       | 21         | [0, 700, 350, 500, 400, 480, 450, 15500, 1200,...] |                   | <i>int</i>              |
| <b><i>3SsnPorch</i></b>     | 20         | [0, 320, 407, 130, 180, 168, 140, 508, 238, 24...  |                   | <i>int</i>              |
| <b><i>Exterior2nd</i></b>   | 16         | [VinylSd, MetalSd, Wd Shng, HdBoard, Plywood, ...] |                   | <i>cat</i>              |
| <b><i>MSSubClass</i></b>    | 15         | [60, 20, 70, 50, 190, 45, 90, 120, 30, 85, 80,...] |                   | <i>cat</i>              |
| <b><i>Exterior1st</i></b>   | 15         | [VinylSd, MetalSd, Wd Sdng, HdBoard, BrkFace, ...] |                   | <i>cat</i>              |
| <b><i>MoSold</i></b>        | 12         | [2, 5, 9, 12, 10, 8, 11, 4, 1, 7, 3, 6]            |                   | <i>cat</i>              |
| <b><i>TotRmsAbvGrd</i></b>  | 12         | [8, 6, 7, 9, 5, 11, 4, 10, 12, 3, 2, 14]           |                   | <i>cat</i>              |
| <b><i>OverallQual</i></b>   | 10         | [7, 6, 8, 5, 9, 4, 10, 3, 1, 2]                    |                   | <i>cat</i>              |
| <b><i>SaleType</i></b>      | 9          | [WD, New, COD, ConLD, ConLI, CWD, ConLw, Con, ...] |                   | <i>cat</i>              |
| <b><i>OverallCond</i></b>   | 9          | [5, 8, 6, 7, 4, 2, 3, 9, 1]                        |                   | <i>cat</i>              |
| <b><i>Condition1</i></b>    | 9          | [Norm, Feedr, PosN, Artery, RRAe, RRNn, RRAn, ...] |                   | <i>cat</i>              |

|                      | <b>Len</b> | <b>Categories</b>                                   | <b>Correct_datatype</b> |
|----------------------|------------|---|-------------------------|
| <b>BedroomAbvGr</b>  | 8          | [3, 4, 1, 2, 0, 5, 6, 8]                            | cat                     |
| <b>PoolArea</b>      | 8          | [0, 512, 648, 576, 555, 480, 519, 738]              | int                     |
| <b>RoofMatl</b>      | 8          | [CompShg, WdShngl, Metal, WdShake, Membran, Ta...]  | cat                     |
| <b>Condition2</b>    | 8          | [Norm, Artery, RRNn, Feedr, PosN, PosA, RRAn, ...]  | cat                     |
| <b>HouseStyle</b>    | 8          | [2Story, 1Story, 1.5Fin, 1.5Unf, SFoyer, SLvl,...]  | cat                     |
| <b>BsmtFinType2</b>  | 7          | [Unf, BLQ, nan, ALQ, Rec, LwQ, GLQ]                 | cat                     |
| <b>GarageType</b>    | 7          | [Attchd, Detachd, BuiltIn, CarPort, nan, Basmen...] | cat                     |
| <b>BsmtFinType1</b>  | 7          | [GLQ, ALQ, Unf, Rec, BLQ, nan, LwQ]                 | cat                     |
| <b>Functional</b>    | 7          | [Typ, Min1, Maj1, Min2, Mod, Maj2, Sev]             | cat                     |
| <b>Electrical</b>    | 6          | [SBrkr, FuseF, FuseA, FuseP, Mix, nan]              | cat                     |
| <b>RoofStyle</b>     | 6          | [Gable, Hip, Gambrel, Mansard, Flat, Shed]          | cat                     |
| <b>Heating</b>       | 6          | [GasA, GasW, Grav, Wall, OthW, Floor]               | cat                     |
| <b>GarageQual</b>    | 6          | [TA, Fa, Gd, nan, Ex, Po]                           | cat                     |
| <b>Foundation</b>    | 6          | [PConc, CBlock, BrkTil, Wood, Slab, Stone]          | cat                     |
| <b>SaleCondition</b> | 6          | [Normal, Abnrmal, Partial, AdjLand, Alloca, Fa...]  | cat                     |
| <b>FireplaceQu</b>   | 6          | [nan, TA, Gd, Fa, Ex, Po]                           | cat                     |
| <b>GarageCond</b>    | 6          | [TA, Fa, nan, Gd, Po, Ex]                           | cat                     |
| <b>MiscFeature</b>   | 5          | [nan, Shed, Gar2, Othr, TenC]                       | cat                     |
| <b>Fence</b>         | 5          | [nan, MnPrv, GdWo, GdPrv, MnWw]                     | cat                     |
| <b>GarageCars</b>    | 5          | [2, 3, 1, 0, 4]                                     | cat                     |
| <b>YrSold</b>        | 5          | [2008, 2007, 2006, 2009, 2010]                      | cat                     |
| <b>HeatingQC</b>     | 5          | [Ex, Gd, TA, Fa, Po]                                | cat                     |
| <b>LotConfig</b>     | 5          | [Inside, FR2, Corner, CulDSac, FR3]                 | cat                     |
| <b>BldgType</b>      | 5          | [1Fam, 2fmCon, Duplex, TwnhsE, Twnhs]               | cat                     |
| <b>MasVnrType</b>    | 5          | [BrkFace, None, Stone, BrkCmn, nan]                 | cat                     |
| <b>ExterCond</b>     | 5          | [TA, Gd, Fa, Po, Ex]                                | cat                     |
| <b>MSZoning</b>      | 5          | [RL, RM, C (all), FV, RH]                           | cat                     |
| <b>BsmtQual</b>      | 5          | [Gd, TA, Ex, nan, Fa]                               | cat                     |
| <b>BsmtCond</b>      | 5          | [TA, Gd, nan, Fa, Po]                               | cat                     |
| <b>BsmtExposure</b>  | 5          | [No, Gd, Mn, Av, nan]                               | cat                     |
| <b>KitchenQual</b>   | 4          | [Gd, TA, Ex, Fa]                                    | cat                     |
| <b>Fireplaces</b>    | 4          | [0, 1, 2, 3]  | cat                     |
| <b>ExterQual</b>     | 4          | [Gd, TA, Ex, Fa]                                    | cat                     |
| <b>PoolQC</b>        | 4          | [nan, Ex, Fa, Gd]                                   | cat                     |

|                            | <i>Len</i> | <i>Categories</i>    | <i>Correct_datatype</i> |
|----------------------------|------------|----------------------|-------------------------|
| <b><i>LotShape</i></b>     | 4          | [Reg, IR1, IR2, IR3] | <i>cat</i>              |
| <b><i>LandContour</i></b>  | 4          | [Lvl, Bnk, Low, HLS] | <i>cat</i>              |
| <b><i>BsmtFullBath</i></b> | 4          | [1, 0, 2, 3]         | <i>cat</i>              |
| <b><i>GarageFinish</i></b> | 4          | [RFn, Unf, Fin, nan] | <i>cat</i>              |
| <b><i>FullBath</i></b>     | 4          | [2, 1, 3, 0]         | <i>cat</i>              |
| <b><i>KitchenAbvGr</i></b> | 4          | [1, 2, 3, 0]         | <i>cat</i>              |
| <b><i>PavedDrive</i></b>   | 3          | [Y, N, P]            | <i>cat</i>              |
| <b><i>BsmtHalfBath</i></b> | 3          | [0, 1, 2]            | <i>cat</i>              |
| <b><i>HalfBath</i></b>     | 3          | [1, 0, 2]            | <i>cat</i>              |
| <b><i>LandSlope</i></b>    | 3          | [Gtl, Mod, Sev]      | <i>cat</i>              |
| <b><i>Alley</i></b>        | 3          | [nan, Grvl, Pave]    | <i>cat</i>              |
| <b><i>Street</i></b>       | 2          | [Pave, Grvl]         | <i>cat</i>              |
| <b><i>CentralAir</i></b>   | 2          | [Y, N]               | <i>cat</i>              |
| <b><i>Utilities</i></b>    | 2          | [AllPub, NoSeWa]     | <i>cat</i>              |

Each value in the column Categories from dataframe df\_unique corresponds to:

| <i>int</i> | <i>float</i> | <i>cat</i> |
|------------|--------------|------------|
| Integer    | Float        | Category   |

Proceed to converted the columns into its correspond type showed in the datafarame df\_unique.

```
In [11]: for x in df_unique.reset_index().values:
    if x[3]=='int':
        df_train[x[0]]=df_train[x[0]].astype('int64')
    elif x[3]=='float':
        df_train[x[0]]=df_train[x[0]].astype('float64')
    elif x[3]=='cat':
        df_train[x[0]]=df_train[x[0]].astype('category')
```

```
In [12]: df_train.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1460 entries, 0 to 1459
Data columns (total 81 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   Id                1460 non-null    int64  
 1   MSSubClass         1460 non-null    category
 2   MSZoning          1460 non-null    category
 3   LotFrontage        1201 non-null    float64 
 4   LotArea            1460 non-null    int64  
 5   Street             1460 non-null    category
 6   Alley              91 non-null     category
 7   LotShape            1460 non-null    category
 8   LandContour        1460 non-null    category
 9   Utilities           1460 non-null    category
 10  LotConfig           1460 non-null    category
 11  LandSlope           1460 non-null    category
 12  Neighborhood        1460 non-null    category
 13  Condition1          1460 non-null    category
 14  Condition2          1460 non-null    category
 15  BldgType            1460 non-null    category
 16  HouseStyle          1460 non-null    category
 17  OverallQual         1460 non-null    category
 18  OverallCond         1460 non-null    category
 19  YearBuilt            1460 non-null    category
 20  YearRemodAdd        1460 non-null    category
 21  RoofStyle            1460 non-null    category
 22  RoofMatl             1460 non-null    category
 23  Exterior1st          1460 non-null    category
 24  Exterior2nd          1460 non-null    category
 25  MasVnrType           1452 non-null    category
 26  MasVnrArea           1452 non-null    float64 
 27  ExterQual            1460 non-null    category
 28  ExterCond            1460 non-null    category
 29  Foundation           1460 non-null    category
 30  BsmtQual             1423 non-null    category
 31  BsmtCond              1423 non-null    category
 32  BsmtExposure          1422 non-null    category
 33  BsmtFinType1          1423 non-null    category
 34  BsmtFinSF1            1460 non-null    int64  
 35  BsmtFinType2          1422 non-null    category
 36  BsmtFinSF2            1460 non-null    int64  
 37  BsmtUnfSF             1460 non-null    int64  
 38  TotalBsmtSF           1460 non-null    int64  
 39  Heating               1460 non-null    category
 40  HeatingQC              1460 non-null    category
 41  CentralAir             1460 non-null    category
 42  Electrical             1459 non-null    category
 43  1stFlrSF              1460 non-null    int64  
 44  2ndFlrSF              1460 non-null    int64  
 45  LowQualFinSF           1460 non-null    int64  
 46  GrLivArea              1460 non-null    int64  
 47  BsmtFullBath           1460 non-null    category
 48  BsmtHalfBath           1460 non-null    category
 49  FullBath               1460 non-null    category
 50  HalfBath                1460 non-null    category
 51  BedroomAbvGr            1460 non-null    category
 52  KitchenAbvGr            1460 non-null    category
 53  KitchenQual             1460 non-null    category
 54  TotRmsAbvGrd            1460 non-null    category

```

```

55  Functional      1460 non-null  category
56  Fireplaces       1460 non-null  category
57  FirePlaceQu     770 non-null   category
58  GarageType       1379 non-null  category
59  GarageYrBlt     1379 non-null  category
60  GarageFinish     1379 non-null  category
61  GarageCars        1460 non-null  category
62  GarageArea        1460 non-null  int64
63  GarageQual       1379 non-null  category
64  GarageCond       1379 non-null  category
65  PavedDrive       1460 non-null  category
66  WoodDeckSF        1460 non-null  int64
67  OpenPorchSF       1460 non-null  int64
68  EnclosedPorch     1460 non-null  int64
69  3SsnPorch        1460 non-null  int64
70  ScreenPorch       1460 non-null  int64
71  PoolArea          1460 non-null  int64
72  PoolQC            7 non-null    category
73  Fence              281 non-null  category
74  MiscFeature       54 non-null   category
75  MiscVal           1460 non-null  int64
76  MoSold            1460 non-null  category
77  YrSold             1460 non-null  category
78  SaleType           1460 non-null  category
79  SaleCondition      1460 non-null  category
80  SalePrice          1460 non-null  int64
dtypes: category(60), float64(2), int64(19)
memory usage: 350.1 KB

```

We got **60 Categorical** variables , **18 Numerical** Variables (excluding *Id*) and **2 float** Variables.

```
In [13]: categories=list(filter(lambda x: x if df_train[x].dtypes=='category' else False,df_train.dtypes))
floats=list(filter(lambda x: x if df_train[x].dtypes=='float64' else False,df_train.dtypes))
integers=list(filter(lambda x: x if df_train[x].dtypes=='int64' else False,df_train.dtypes))
integers.remove('Id')
```

Once the conversion is succeeded, let's use `df.describe()`

```
In [14]: df_train.describe()
```

Out[14]:

|                     | <b><i>Id</i></b> | <b><i>LotFrontage</i></b> | <b><i>LotArea</i></b> | <b><i>MasVnrArea</i></b> | <b><i>BsmtFinSF1</i></b> | <b><i>BsmtFinSF2</i></b> | <b><i>Bsr</i></b> |
|---------------------|------------------|---------------------------|-----------------------|--------------------------|--------------------------|--------------------------|-------------------|
| <b><i>count</i></b> | 1460.000000      | 1201.000000               | 1460.000000           | 1452.000000              | 1460.000000              | 1460.000000              | 1460.000000       |
| <b><i>mean</i></b>  | 730.500000       | 70.049958                 | 10516.828082          | 103.685262               | 443.639726               | 46.549315                | 567               |
| <b><i>std</i></b>   | 421.610009       | 24.284752                 | 9981.264932           | 181.066207               | 456.098091               | 161.319273               | 441               |
| <b><i>min</i></b>   | 1.000000         | 21.000000                 | 1300.000000           | 0.000000                 | 0.000000                 | 0.000000                 | 0                 |
| <b><i>25%</i></b>   | 365.750000       | 59.000000                 | 7553.500000           | 0.000000                 | 0.000000                 | 0.000000                 | 223               |
| <b><i>50%</i></b>   | 730.500000       | 69.000000                 | 9478.500000           | 0.000000                 | 383.500000               | 0.000000                 | 477               |
| <b><i>75%</i></b>   | 1095.250000      | 80.000000                 | 11601.500000          | 166.000000               | 712.250000               | 0.000000                 | 808               |
| <b><i>max</i></b>   | 1460.000000      | 313.000000                | 215245.000000         | 1600.000000              | 5644.000000              | 1474.000000              | 2336              |

◀ ▶

## Missing Data

Actual data with null values are:

In [15]:

```
for x in df_train.columns:
    if (df_train[x].isna().sum() > 0):
        print(' {}: {} records are null out of {}. This means {:.2f}% of all rec
```

*LotFrontage: 259 records are null out of 1460. This means 17.74% of all records are missing*

*Alley: 1369 records are null out of 1460. This means 93.77% of all records are missing*

*MasVnrType: 8 records are null out of 1460. This means 0.55% of all records are missing*

*MasVnrArea: 8 records are null out of 1460. This means 0.55% of all records are missing*

*BsmtQual: 37 records are null out of 1460. This means 2.53% of all records are missing*

*BsmtCond: 37 records are null out of 1460. This means 2.53% of all records are missing*

*BsmtExposure: 38 records are null out of 1460. This means 2.60% of all records are missing*

*BsmtFinType1: 37 records are null out of 1460. This means 2.53% of all records are missing*

*BsmtFinType2: 38 records are null out of 1460. This means 2.60% of all records are missing*

*Electrical: 1 records are null out of 1460. This means 0.07% of all records are missing*

*FireplaceQu: 690 records are null out of 1460. This means 47.26% of all records are missing*

*GarageType: 81 records are null out of 1460. This means 5.55% of all records are missing*

*GarageYrBlt: 81 records are null out of 1460. This means 5.55% of all records are missing*

*GarageFinish: 81 records are null out of 1460. This means 5.55% of all records are missing*

*GarageQual: 81 records are null out of 1460. This means 5.55% of all records are missing*

*GarageCond: 81 records are null out of 1460. This means 5.55% of all records are missing*

*PoolQC: 1453 records are null out of 1460. This means 99.52% of all records are missing*

*Fence: 1179 records are null out of 1460. This means 80.75% of all records are missing*

*MiscFeature: 1406 records are null out of 1460. This means 96.30% of all records are missing*

To see a general view of proportion of null values all over the data.

The next graph helps to analyze better this event

```
In [16]: missing={}
for i in df_train:
    if df_train[i].isna().sum()>0:
        missing.update({i:df_train[i].isna().sum()})
```

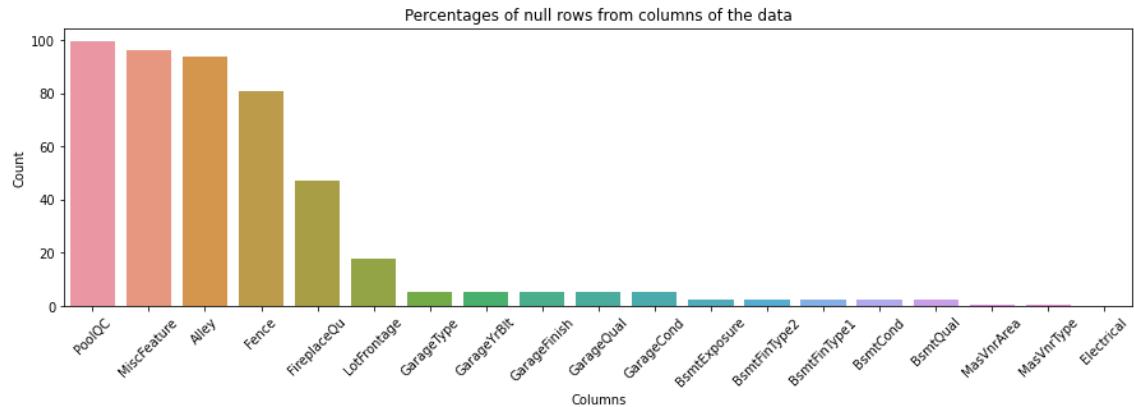
```
In [17]: df_missing=pd.DataFrame(missing,index=['Count']).transpose().sort_values(by='Count', ascending=False)
df_missing.columns=['missing']
for x in df_missing.index:
    df_missing.loc[x,'missing_perc']=df_missing.loc[x,'missing']/len(df_train)
```

```
In [18]: import seaborn as sns
import matplotlib.pyplot as plt
plt.figure(figsize=(15,4))
sns.barplot(x=df_missing.index,y=df_missing.missing_perc)
plt.title('Percentages of null rows from columns of the data')
plt.xlabel('Columns')
```

```

plt.ylabel('Count')
plt.xticks(rotation=45)
#plt.plot([0, 18], [10, 10], 'g', lw=2)
plt.show()

```



**From this, it can be seen:**

About **100%** of houses do not have Pool, Miscellaneous feature and No alley access

More than **80%** of records do not have Fence

About **50%** of records do not have fireplaces

This means that the empty values not always mean missing values but tell **another story**.

**However, there are others variables which we cannot say the same.**

**For example:**

- **LotFrontage: Null values in the linear feet of street connected to property. It might be data quality lack.**
- **MSZoning: Null values in the general zoning classification of the sale. It might be data quality lack too.**
- **So on...**

**So It is important to proceed a data cleaning process.**

## Data Cleaning

---

**In this section, some attributes will be filled due to its values are important for the later analysis**

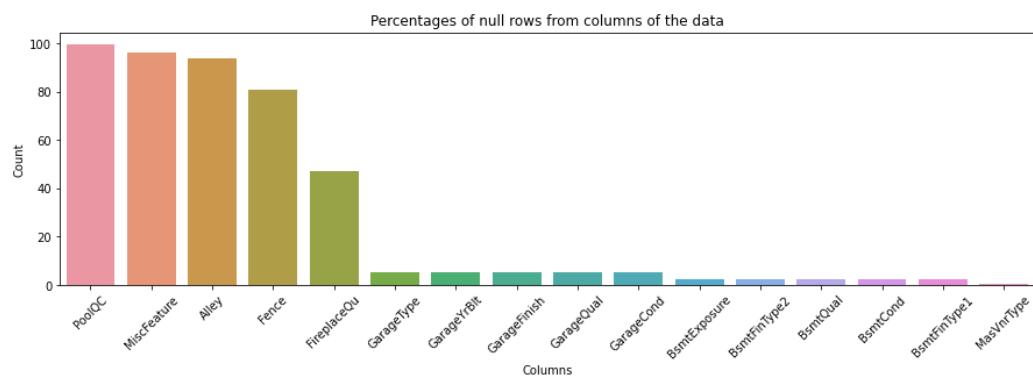
```
In [19]: # numerical
df_train.LotFrontage=df_train.groupby("Neighborhood")["LotFrontage"].transform()
df_train.MasVnrArea=df_train.MasVnrArea.fillna(0) # Those values in Nan for categorical
df_train['MSZoning']=df_train.groupby("Neighborhood")["MSZoning"].transform()
df_train['Electrical']=df_train.groupby("Neighborhood")["Electrical"].transform()
```

**Since the previous code was only substituted the columns: LotFrontage, MSZoning, MasVnrArea and Electrical, the rest of the attributes with null values can be substituted by None because this value makes sense in the categories.**

```
In [20]: missing={}
for i in df_train:
    if df_train[i].isna().sum()>0:
        missing.update({i:df_train[i].isna().sum()})
```

```
In [21]: df_missing=pd.DataFrame(missing,index=['Count']).transpose().sort_values('Count', ascending=False)
df_missing.columns=['missing']
for x in df_missing.index:
    df_missing.loc[x,'missing_perc']=df_missing.loc[x,'missing']/len(df_train)
```

```
In [22]: import seaborn as sns
import matplotlib.pyplot as plt
plt.figure(figsize=(15,4))
sns.barplot(x=df_missing.index,y=df_missing.missing_perc)
plt.title('Percentages of null rows from columns of the data')
plt.xlabel('Columns')
plt.ylabel('Count')
plt.xticks(rotation=45)
#plt.plot([0, 18], [10, 10], 'g', lw=2)
plt.show()
```



**Now, it left 16 attributes with NaN**

```
df_train['PoolQC']=df_train['PoolQC'].cat.add_categories(['None'])
```

In [23]:

```
attri_null=pd.DataFrame(df_train.isnull().sum()[df_train.isnull().sum()
for i in attri_null.index:
    print(i)
    if 'None' not in df_train[i].cat.categories:
        df_train[i]=df_train[i].cat.add_categories(['None'])
        df_train[i] = df_train[i].fillna('None')
```

PoolQC  
MiscFeature  
Alley  
Fence  
FireplaceQu  
GarageType  
GarageYrBlt  
GarageFinish  
GarageQual  
GarageCond  
BsmtExposure  
BsmtFinType2  
BsmtQual  
BsmtCond  
BsmtFinType1  
MasVnrType

In [24]:

```
print('There are {} columns with NULL values'.format(len(df_train.isna().sum())))
There are 0 columns with NULL values
```

**Great!. There is not more NULL values in the dataset.**

## Univariable analysis

### Numerical variables

In [25]:

```
def plot_subplot(itera,n,Label,compare=' '):
    row=int(len(itera)/n)+len(itera)%n
    colu=len(itera)-row*n

    if Label=='hist':
        plt.figure(figsize=(16,6*row))
        for i,j in enumerate(itera):
            i+=1
            plt.subplot(row,n,i)
            sns.violinplot(data=df_train,y=j,kde=True,hue='Neighborhood')
            plt.title('Distribution of '+j)
    elif Label=='cross':
        a=0
        b=0
        fig, axes=plt.subplots(row,n,figsize=(16,8*row))
```

```

# print(row,n)
for i,j in enumerate(itera):
    i+=1
    my_crosstab = df_train[j].value_counts(sort=False)/sum(my_crosstab)
    if row==1 or n==1:
        #print('a',a,b)
        ax=my_crosstab.plot(kind='bar', stacked=True, rot=90)
        #ax.set_xlabel(j)
        ax.set_ylabel('% Count')
    else:
        #print(a,b)
        ax=my_crosstab.plot(kind='bar', stacked=True, rot=90)
        ax.set_xlabel(j)
        ax.set_ylabel('% Count')
    if b==n-1:
        b=0
        a+=1
    else:
        b+=1

```

In [26]: `df_train.describe().round(1)`

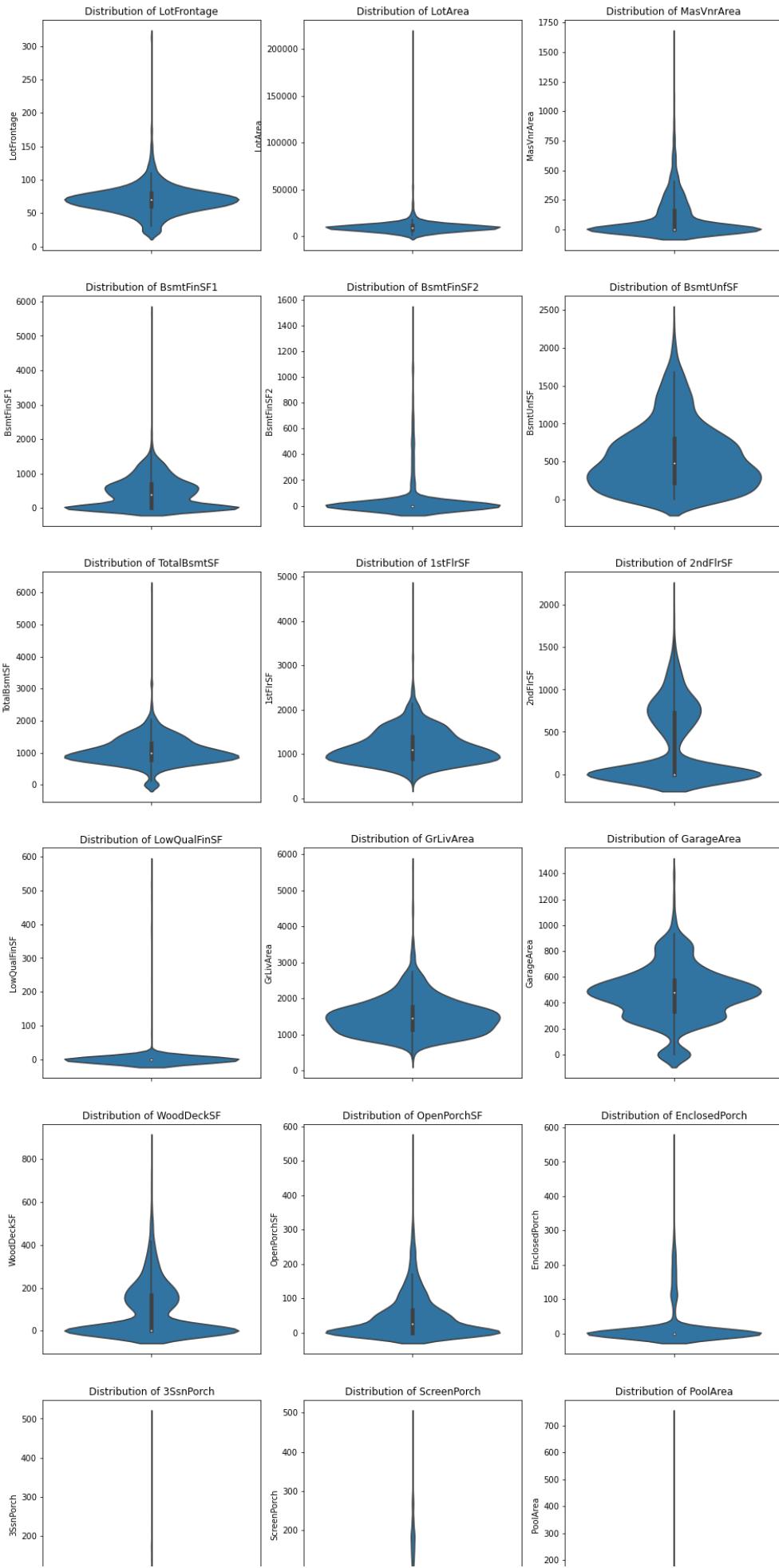
Out[26]:

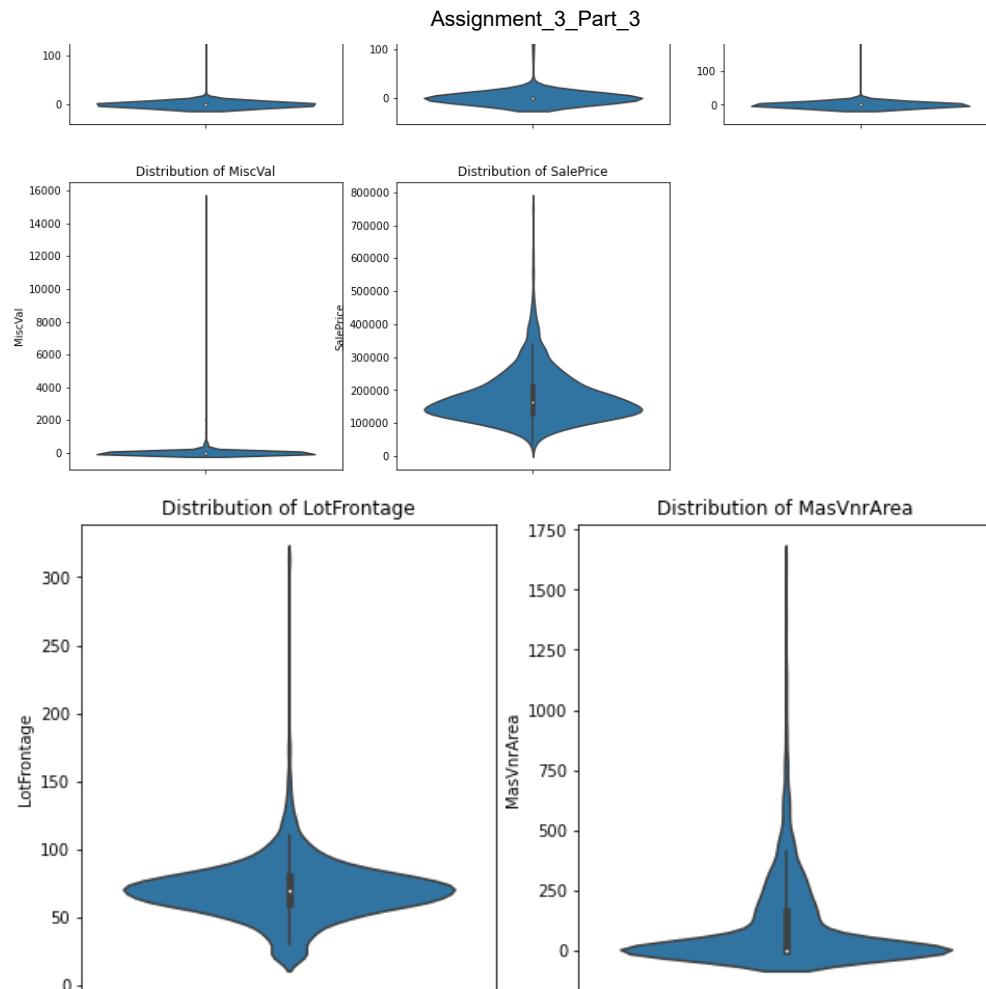
|                            | <u><b><i>Id</i></b></u>     | <u><b><i>LotFrontage</i></b></u> | <u><b><i>LotArea</i></b></u>  | <u><b><i>MasVnrArea</i></b></u> | <u><b><i>BsmtFinSF1</i></b></u> | <u><b><i>BsmtFinSF2</i></b></u> | <u><b><i>BsmtFinTottsf</i></b></u> |
|----------------------------|-----------------------------|----------------------------------|-------------------------------|---------------------------------|---------------------------------|---------------------------------|------------------------------------|
| <u><b><i>count</i></b></u> | <u><b><i>1460.0</i></b></u> | <u><b><i>1460.0</i></b></u>      | <u><b><i>1460.0</i></b></u>   | <u><b><i>1460.0</i></b></u>     | <u><b><i>1460.0</i></b></u>     | <u><b><i>1460.0</i></b></u>     | <u><b><i>1460.0</i></b></u>        |
| <u><b><i>mean</i></b></u>  | <u><b><i>730.5</i></b></u>  | <u><b><i>70.2</i></b></u>        | <u><b><i>10516.8</i></b></u>  | <u><b><i>103.1</i></b></u>      | <u><b><i>443.6</i></b></u>      | <u><b><i>46.5</i></b></u>       |                                    |
| <u><b><i>std</i></b></u>   | <u><b><i>421.6</i></b></u>  | <u><b><i>22.4</i></b></u>        | <u><b><i>9981.3</i></b></u>   | <u><b><i>180.7</i></b></u>      | <u><b><i>456.1</i></b></u>      | <u><b><i>161.3</i></b></u>      |                                    |
| <u><b><i>min</i></b></u>   | <u><b><i>1.0</i></b></u>    | <u><b><i>21.0</i></b></u>        | <u><b><i>1300.0</i></b></u>   | <u><b><i>0.0</i></b></u>        | <u><b><i>0.0</i></b></u>        | <u><b><i>0.0</i></b></u>        |                                    |
| <u><b><i>25%</i></b></u>   | <u><b><i>365.8</i></b></u>  | <u><b><i>60.0</i></b></u>        | <u><b><i>7553.5</i></b></u>   | <u><b><i>0.0</i></b></u>        | <u><b><i>0.0</i></b></u>        | <u><b><i>0.0</i></b></u>        |                                    |
| <u><b><i>50%</i></b></u>   | <u><b><i>730.5</i></b></u>  | <u><b><i>70.0</i></b></u>        | <u><b><i>9478.5</i></b></u>   | <u><b><i>0.0</i></b></u>        | <u><b><i>383.5</i></b></u>      | <u><b><i>0.0</i></b></u>        |                                    |
| <u><b><i>75%</i></b></u>   | <u><b><i>1095.2</i></b></u> | <u><b><i>80.0</i></b></u>        | <u><b><i>11601.5</i></b></u>  | <u><b><i>164.2</i></b></u>      | <u><b><i>712.2</i></b></u>      | <u><b><i>0.0</i></b></u>        |                                    |
| <u><b><i>max</i></b></u>   | <u><b><i>1460.0</i></b></u> | <u><b><i>313.0</i></b></u>       | <u><b><i>215245.0</i></b></u> | <u><b><i>1600.0</i></b></u>     | <u><b><i>5644.0</i></b></u>     | <u><b><i>1474.0</i></b></u>     |                                    |

In [27]: `integers=list(df_train.describe().round(1).columns.drop('Id'))`

In [28]: `plot_subplot(integers,3,'hist')`  
`plot_subplot(floats,3,'hist')`

## Assignment\_3\_Part\_3

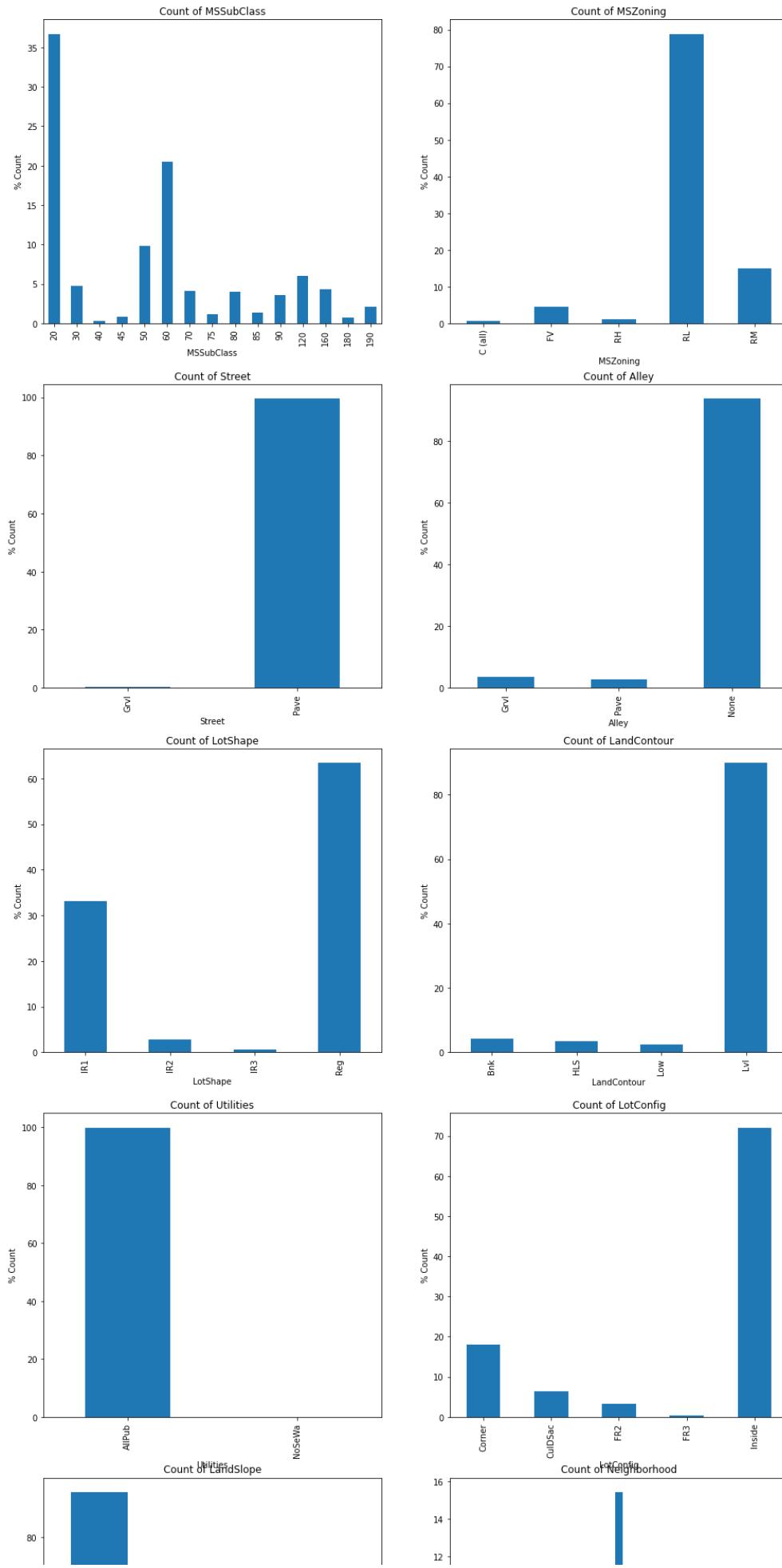




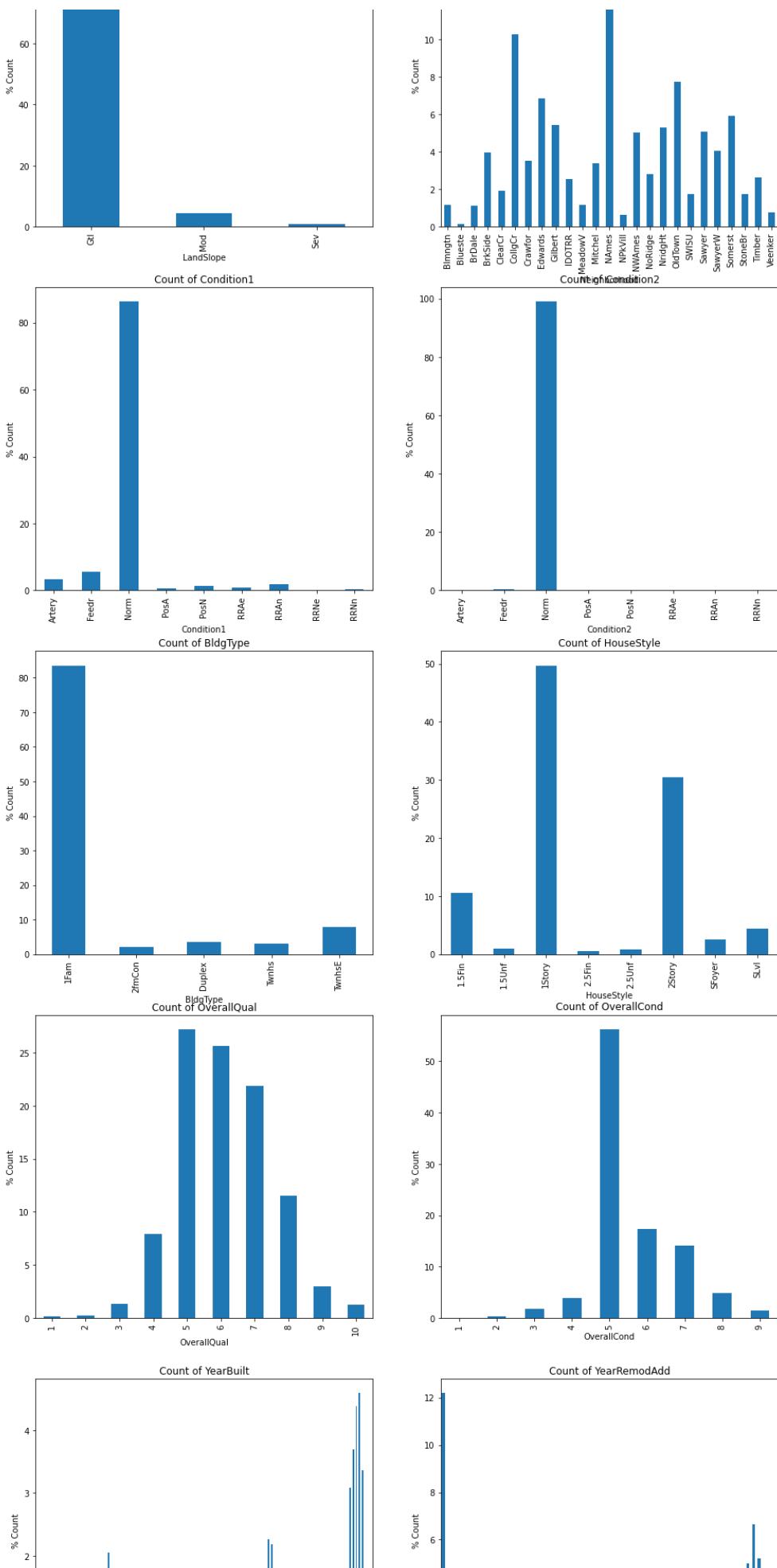
## Categorical variables

In [29]: `plot_subplot(categories,2,'cross','GarageQual')`

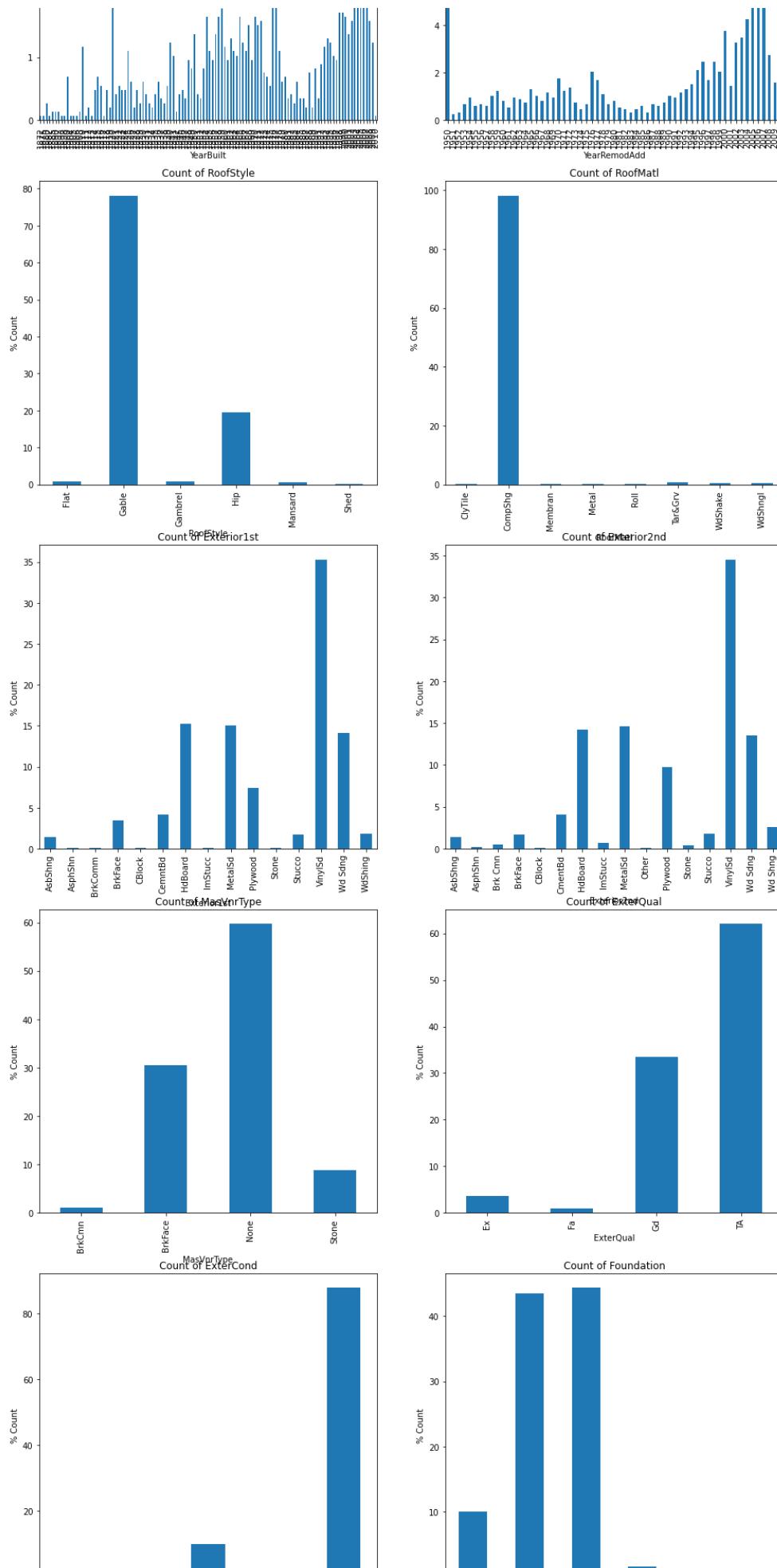
## Assignment\_3\_Part\_3



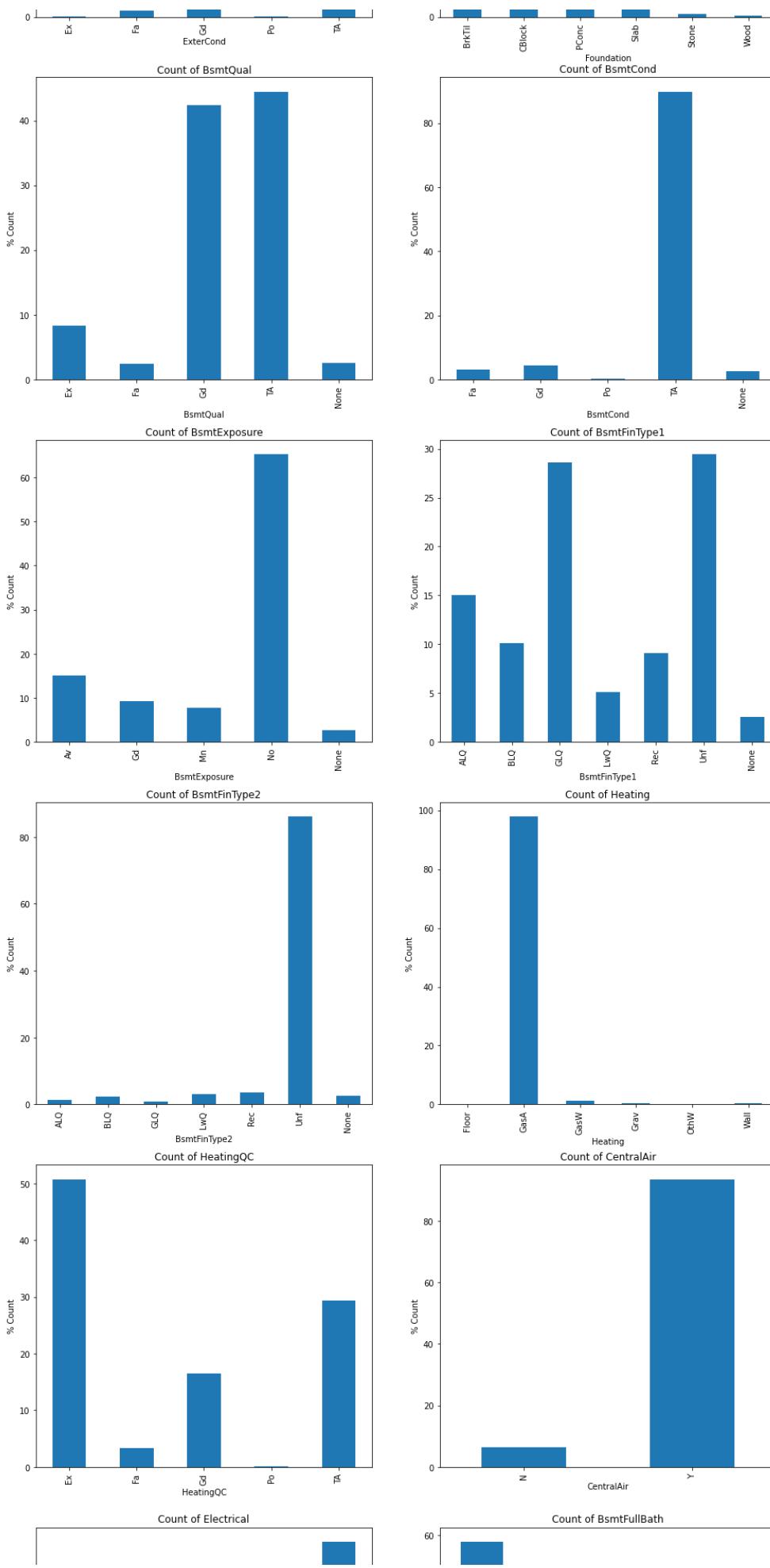
## Assignment\_3\_Part\_3



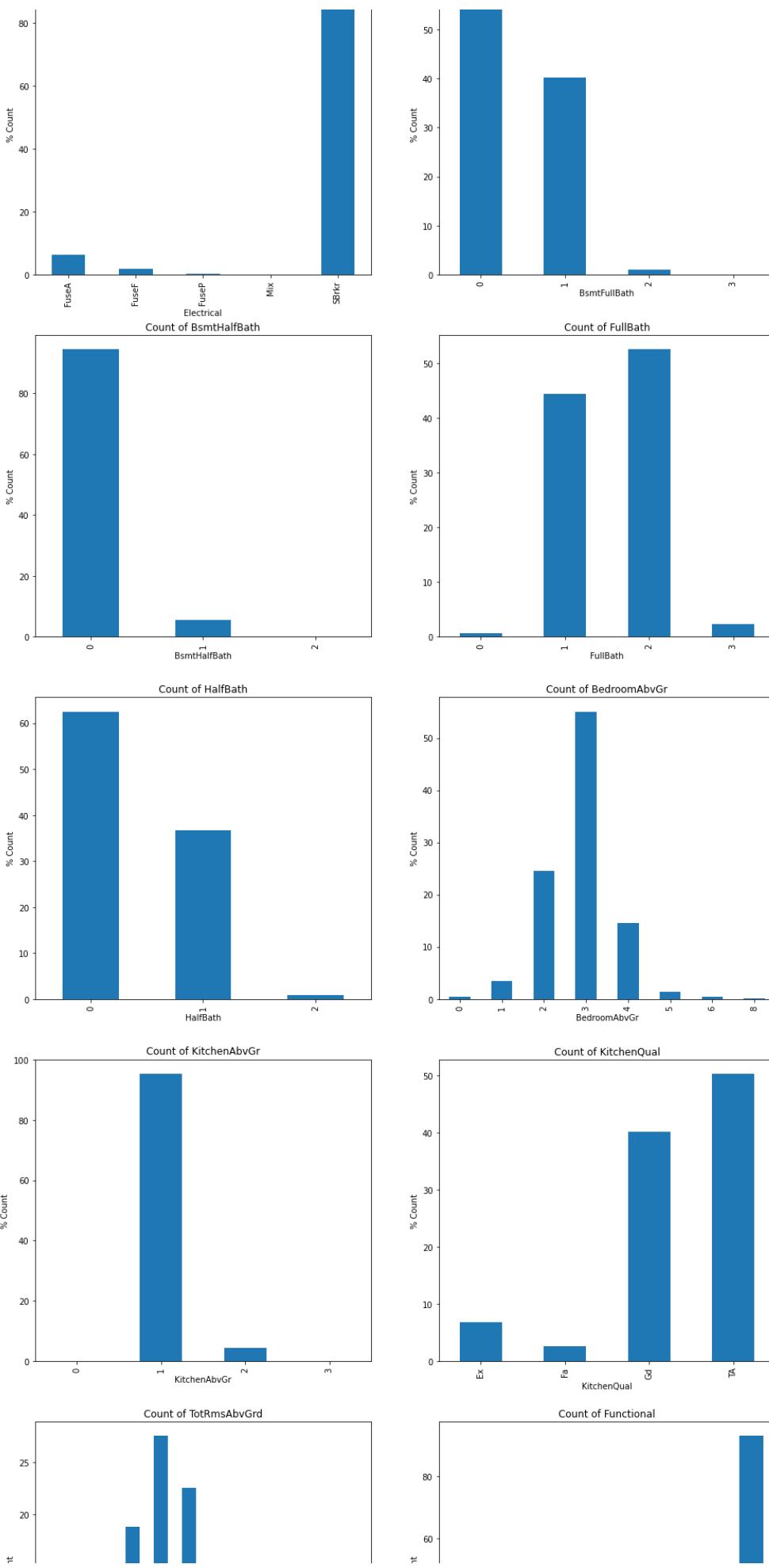
## Assignment\_3\_Part\_3



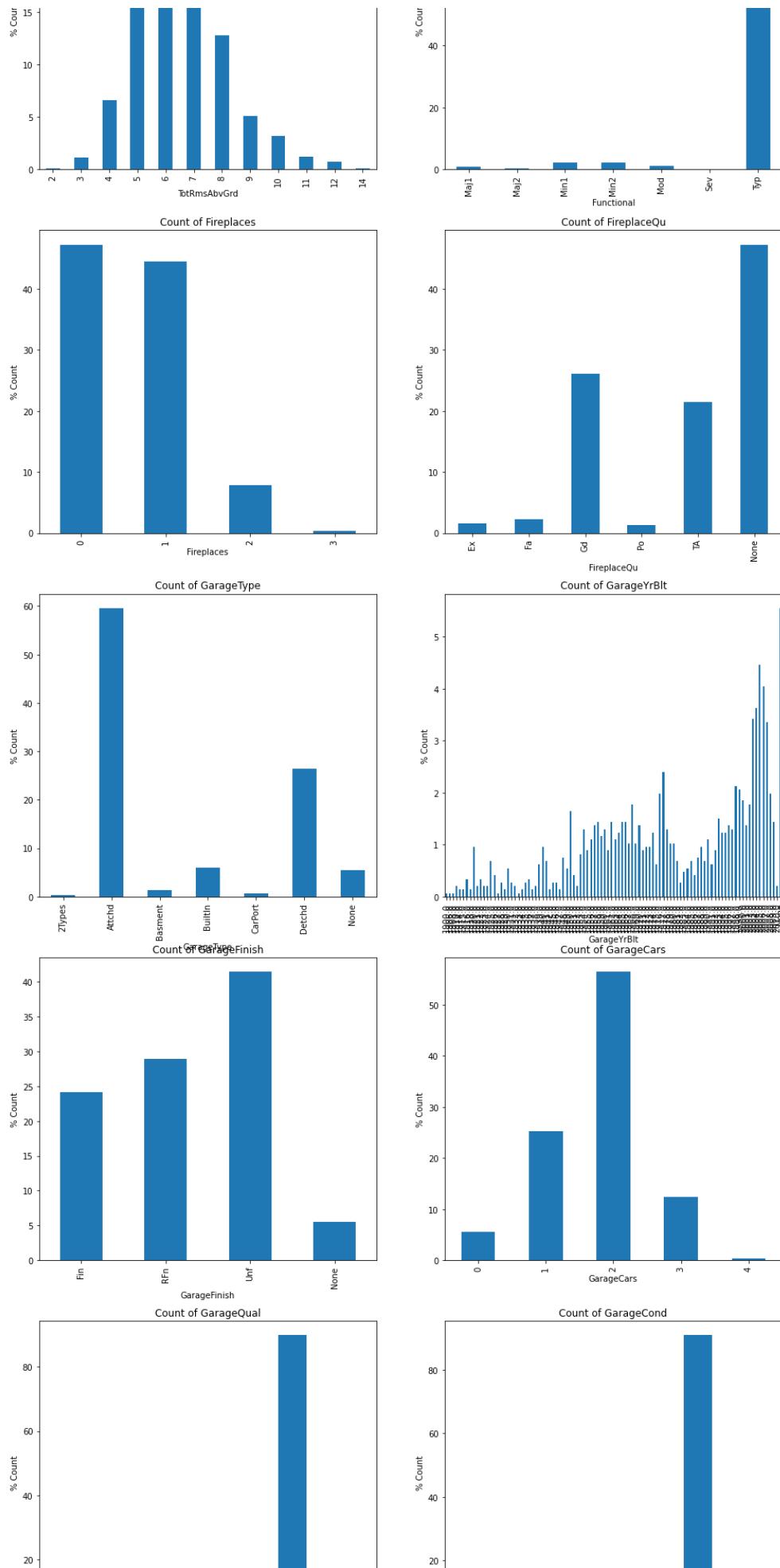
## Assignment\_3\_Part\_3



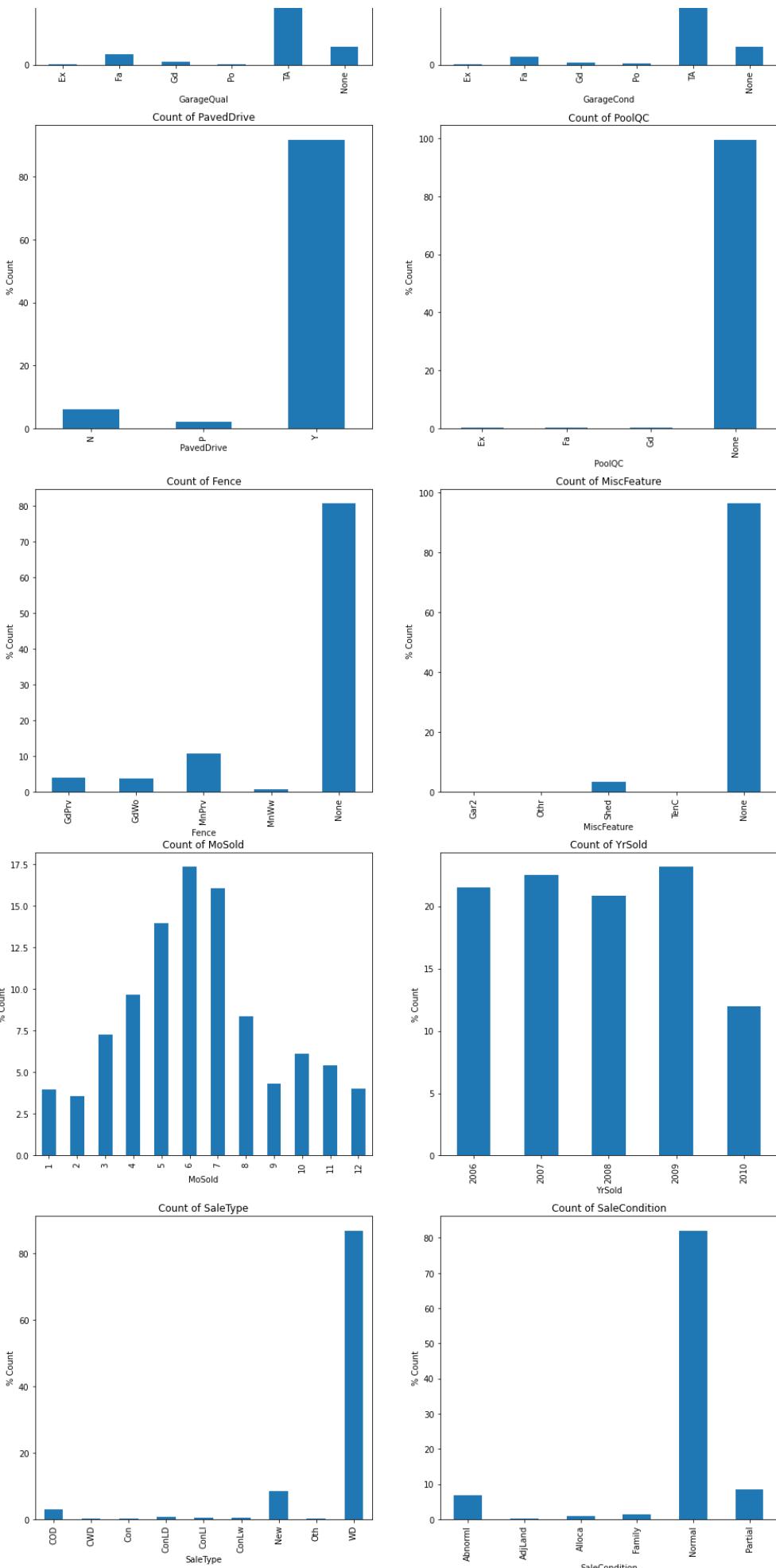
## Assignment\_3\_Part\_3



## Assignment\_3\_Part\_3



## Assignment\_3\_Part\_3

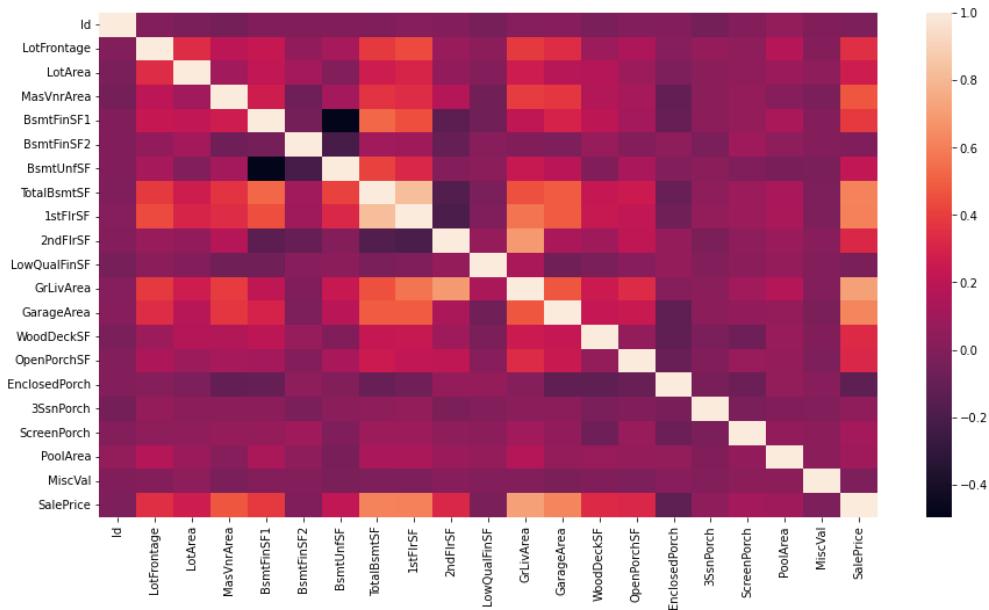


# Bivariable analysis

## Numerical variables

In [30]: `plt.figure(figsize=(15,8))  
sns.heatmap(df_train.corr())`

Out[30]: `<AxesSubplot:>`



## Categorical variables

### Hypothesis 1

This analysis will be answered within some hypothesis proposed.

**Is the Exterior covering on house associated with lot configuration ?**

#### LotConfig: Lot configuration

| <u>Category</u> | <u>Meaning</u>                         |
|-----------------|--|
| <u>Inside</u>   | <u>Inside lot</u>                      |
| <u>Corner</u>   | <u>Corner lot</u>                      |
| <u>CulDSac</u>  | <u>Cul-de-sac</u>                      |
| <u>FR2</u>      | <u>Frontage on 2 sides of property</u> |

| <u>Category</u> | <u>Meaning</u>                         |
|-----------------|--|
| <u>FR3</u>      | <u>Frontage on 3 sides of property</u> |

Exterior1st: Exterior covering on house

| <u>Category</u> | <u>Meaning</u>           |
|-----------------|--------------------------|
| <u>AsbShng</u>  | <u>Asbestos Shingles</u> |
| <u>AsphShn</u>  | <u>Asphalt Shingles</u>  |
| <u>BrkComm</u>  | <u>Brick Common</u>      |
| <u>BrkFace</u>  | <u>Brick Face</u>        |
| <u>CBlock</u>   | <u>Cinder Block</u>      |
| <u>CemntBd</u>  | <u>Cement Board</u>      |
| <u>HdBoard</u>  | <u>Hard Board</u>        |
| <u>ImStucc</u>  | <u>Imitation Stucco</u>  |
| <u>MetalSd</u>  | <u>Metal Siding</u>      |
| <u>Other</u>    | <u>Other</u>             |
| <u>Plywood</u>  | <u>Plywood</u>           |
| <u>PreCast</u>  | <u>PreCast</u>           |
| <u>Stone</u>    | <u>Stone</u>             |
| <u>Stucco</u>   | <u>Stucco</u>            |
| <u>VinylSd</u>  | <u>Vinyl Siding</u>      |
| <u>Wd Sdng</u>  | <u>Wood Siding</u>       |
| <u>WdShing</u>  | <u>Wood Shingles</u>     |

In [31]:

```
def plot_subplot2(itera,n,label,compare=' '):
    row=int(len(itera)/n)+len(itera)%n
    colu=len(itera)-row*n

    if label=='hist':
        plt.figure(figsize=(16,6*row))
        for i,j in enumerate(itera):
            i+=1
            plt.subplot(row,n,i)
            sns.violinplot(data=df_train,y=j,kde=True,hue='Neighborhood')
            plt.title('Distribution of '+j)
    elif label=='cross':
        a=0
        b=0
        fig, axes=plt.subplots(row,n,figsize=(16,6*row))
        #print(row,n)
        for i,j in enumerate(itera):
            i+=1
            my_crosstab = pd.crosstab(index=df_train[j],columns=df_train[compare])
            if row==1:
```

```
#print('a',a,b)
my_crosstab.plot(kind='bar', stacked=True, rot=90)

else:
    #print(a,b)
    my_crosstab.plot(kind='bar', stacked=True, rot=90)
    if b==n-1:
        b=0
        a+=1
    else:
        b+=1
```

In [32]:

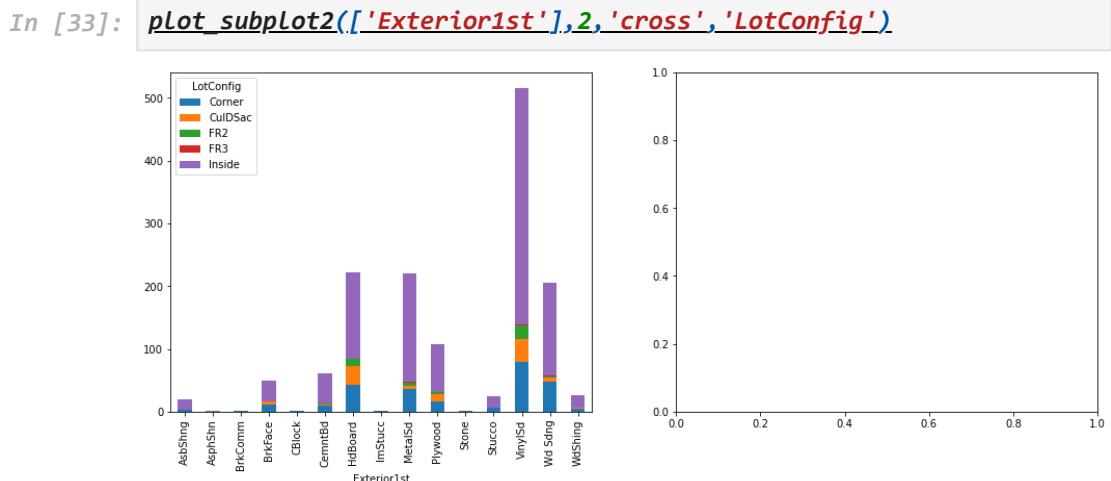
```
def chi_square_mul(itera):
    d={}
    from scipy.stats import chi2_contingency
    for j in itera:
        for i in itera:
            if not i==j:
                my_crosstab = pd.crosstab(index=df_train[i], columns=df_train[j])
                # Chi-square test of independence.
                c, p, dof, expected = chi2_contingency(my_crosstab)
                # Print the p-value
                d.update({(i,j):p})
    return d
```

The hypothesis will be test looks like:

$H_0$ : LotConfig and Exterior1st are independent to each other among all subjects in the population

$H_a$ : LotConfig and Exterior1st are NOT independent to each other among all subjects in the population

## Visual demonstration



## Statistical demonstration

In [34]: `p_value_interest=chi_square_mul([('LotConfig','Exterior1st')))`  
`print('With 95% of confidence, the p-value for LotConfig and Exterior1st is 0.07')`

Taking account this:

- $p(X) < 0.01$  very strong evidence,
- $p(X) \in (0.01, 0.05)$  strong evidence,
- $p(X) \in (0.05, 0.1)$  weak evidence,
- $p(X) > 0.1$  little or no evidence.

This mean as there is a weak evidence to rejected the null hypothesis, it can be concluded that LotConfig and Exterior1st are independent to each other among all population

## Hypothesis 2

Is the roof style of the house associated with Neighborhood ?

RoofStyle: Type of roof

| <u>Category</u> |
|-----------------|
| <u>Flat</u>     |
| <u>Gable</u>    |
| <u>Gambrel</u>  |
| <u>Hip</u>      |
| <u>Mansard</u>  |
| <u>Shed</u>     |

Neighborhood: Physical locations within Ames city limits

| <u>Category</u> | <u>Details</u>             |
|-----------------|----------------------------|
| <u>Blmngtn</u>  | <u>Bloomington Heights</u> |
| <u>Blueste</u>  | <u>Bluestem</u>            |
| <u>BrDale</u>   | <u>Briardale</u>           |
| <u>BrkSide</u>  | <u>Brookside</u>           |

| <u>Category</u> | <u>Details</u>                                   |
|-----------------|--|
| <u>ClearCr</u>  | <u>Clear Creek</u>                               |
| <u>CollgCr</u>  | <u>College Creek</u>                             |
| <u>Crawfor</u>  | <u>Crawford</u>                                  |
| <u>Edwards</u>  | <u>Edwards</u>                                   |
| <u>Gilbert</u>  | <u>Gilbert</u>                                   |
| <u>IDOTRR</u>   | <u>Iowa DOT and Rail Road</u>                    |
| <u>MeadowV</u>  | <u>Meadow Village</u>                            |
| <u>Mitchel</u>  | <u>Mitchell</u>                                  |
| <u>Names</u>    | <u>North Ames</u>                                |
| <u>NoRidge</u>  | <u>Northridge</u>                                |
| <u>NPkVill</u>  | <u>Northpark Villa</u>                           |
| <u>NridgHt</u>  | <u>Northridge Heights</u>                        |
| <u>NWAmes</u>   | <u>Northwest Ames</u>                            |
| <u>OldTown</u>  | <u>Old Town</u>                                  |
| <u>SWISU</u>    | <u>South &amp; West of Iowa State University</u> |
| <u>Sawyer</u>   | <u>Sawyer</u>                                    |
| <u>SawyerW</u>  | <u>Sawyer West</u>                               |
| <u>Somerst</u>  | <u>Somerset</u>                                  |
| <u>StoneBr</u>  | <u>Stone Brook</u>                               |
| <u>Timber</u>   | <u>Timberland</u>                                |
| <u>Veenker</u>  | <u>Veenker</u>                                   |

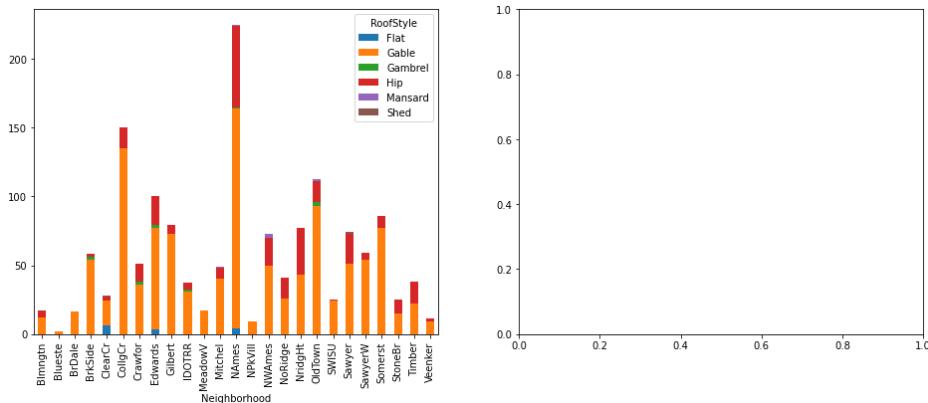
The hypothesis will be test looks like:

$H_0$ : RoofStyle and Neighborhood are independent to each other among all subjects in the population

$H_a$ : RoofStyle and Neighborhood are NOT independent to each other among all subjects in the population

### Visual demonstration

```
In [35]: plot_subPlot2(['Neighborhood'], 2, 'cross', 'RoofStyle')
```



## Statistical demonstration

```
In [36]: p_value_interest=chi_square_mull(('Neighborhood', 'RoofStyle'))
print('With 95% of confidence, the p-value for Neighborhood and RoofStyle is')
p_value_interest
```

With 95% of confidence, the p-value for Neighborhood and RoofStyle is 0.00

Taking account this:

- $p(X) < 0.01$  very strong evidence,
- $p(X) \in (0.01, 0.05)$  strong evidence,
- $p(X) \in (0.05, 0.1)$  weak evidence,
- $p(X) > 0.1$  little or no evidence.

This mean as there is a very strong evidence that null hypothesis can be rejected, it can be concluded that LotConfig and Exterior1st are not significant independent to each other among all population

## Hypothesis 3

Is the type of dwelling involved in the sale associated with the general zoning classification of the sale. ?

MSSubClass: Identifies the type of dwelling involved in the sale.

| <u>Category</u> | <u>Details</u>                             |
|-----------------|--|
| <u>20</u>       | <u>1-STORY 1946 &amp; NEWER ALL STYLES</u> |

| <u>Category</u> | <u>Details</u>   |
|-----------------|--|
| <u>30</u>       | <u>1-STORY 1945 &amp; OLDER</u>                                  |
| <u>40</u>       | <u>1-STORY W/FINISHED ATTIC ALL AGES</u>                         |
| <u>45</u>       | <u>1-1/2 STORY - UNFINISHED ALL AGES</u>                         |
| <u>50</u>       | <u>1-1/2 STORY FINISHED ALL AGES</u>                             |
| <u>60</u>       | <u>2-STORY 1946 &amp; NEWER</u>                                  |
| <u>70</u>       | <u>2-STORY 1945 &amp; OLDER</u>                                  |
| <u>75</u>       | <u>2-1/2 STORY ALL AGES</u>                                      |
| <u>80</u>       | <u>SPLIT OR MULTI-LEVEL</u>                                      |
| <u>85</u>       | <u>SPLIT FOYER</u>   |
| <u>90</u>       | <u>DUPLEX - ALL STYLES AND AGES</u>                              |
| <u>120</u>      | <u>1-STORY PUD (Planned Unit Development) - 1946 &amp; NEWER</u> |
| <u>150</u>      | <u>1-1/2 STORY PUD - ALL AGES</u>                                |
| <u>160</u>      | <u>2-STORY PUD - 1946 &amp; NEWER</u>                            |
| <u>180</u>      | <u>PUD - MULTILEVEL - INCL SPLIT LEV/FOYER</u>                   |
| <u>190</u>      | <u>2 FAMILY CONVERSION - ALL STYLES AND AGES</u>                 |

MSZoning: Identifies the general zoning classification of the sale.

| <u>Category</u> | <u>Meaning</u>                      |
|-----------------|-------------------------------------|
| <u>A</u>        | <u>Agriculture</u>                  |
| <u>C</u>        | <u>Commercial</u>                   |
| <u>FV</u>       | <u>Floating Village Residential</u> |
| <u>I</u>        | <u>Industrial</u>                   |
| <u>RH</u>       | <u>Residential High Density</u>     |
| <u>RL</u>       | <u>Residential Low Density</u>      |
| <u>RP</u>       | <u>Residential Low Density Park</u> |
| <u>RM</u>       | <u>Residential Medium Density</u>   |

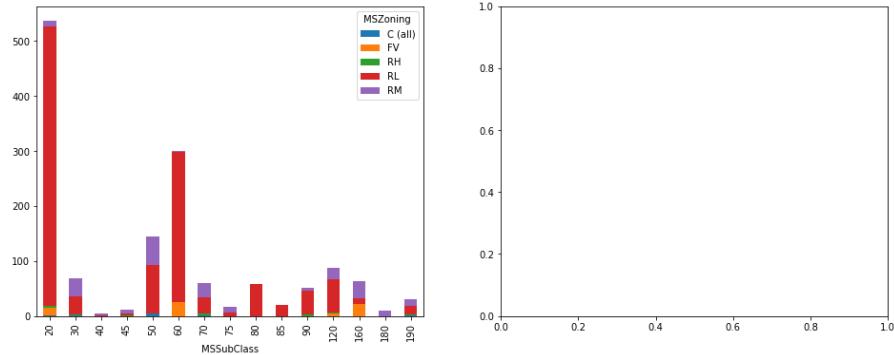
The hypothesis will be test looks like:

$H_0$ : MSSubClass and MSZoning are independent to each other among all subjects in the population

$H_a$ : MSSubClass and MSZoning are NOT independent to each other among all subjects in the population

## Visual demostration

In [37]: `plot_subPlot2(['MSSubClass'], 2, 'cross', 'MSZoning')`



## Statistical demostration

In [38]: `p_value_interest=chi_square_mull(['MSSubClass', 'MSZoning'])`  
`print('With 95% of confidence, the p-value for Neighborhood an`

With 95% of confidence, the p-value for Neighborhood and Roof  
 Style is 0.00

Taking account this:

- $p(X) < 0.01$  very strong evidence.
- $p(X) \in (0.01, 0.05)$  strong evidence.
- $p(X) \in (0.05, 0.1)$  weak evidence,
- $p(X) > 0.1$  little or no evidence.

This mean as there is a very strong evidence that null hypothesis can be rejected, it can be concluded that LotConfig and Exterior1st are not significant independent to each other among all population