

# Exploratory Data Analysis and Supervised Machine Learning (Regression)

## Summary

This work is divided into several sections:

### A) Import the required libraries

### B) Dataset description and understanding data

- Import dataset
- Feature description
- Dataset information
- Analysis of the target (SalePrice)
- Looking for Correlations
- Feature distribution

### C) Data Cleaning

- 1. Handling the Duplicates
- 1. Handling the Missing Values
  - Missing values research
  - Replacement of missing values and deletion of variables
  - Verification of missing values
- 1. Handling the Outliers SF
  - Gr Liv Area
  - Lot Area
  - Low Qual Fin SF
- 1. Feature modification

### D) Exploratory data visualization

- Numerical features
- Categorical features

### E) Hypothesis testing

- Log transformation of sale price
- Central Air
- Land Slope
- MS Zoning
- Summary

### F) Feature Engineering

- 1. Log transforming skew variables
- 1. Basic feature engineering: adding polynomial and interaction terms
  - Feature representation
  - Polynomial Features
  - Interaction terms
- 1. Feature Scaling
- 1. One-hot encoding for dummy variables

### G) Modelisation

- 1. Train test split
- 1. Simple linear regression with GridSearchCV
- 1. Scaling numerical features
- 1. RidgeCV
- 1. LassoCV
- 1. ElasticNetCV
- 1. Result and discussion

## Objective

The main objective of this data analysis is to predict the price of houses based on numerous characteristics. The dataset used is the Ames housing dataset. It examines features of houses sold in Ames (a small city in the state of Iowa in the United States) during the 2006–2010 timeframe.

## Dataset description:

The dataset is composed of 2931 rows and 81 columns. 'SalePrice' is our target or response variable and the rest of the characteristics are our predictor features. We also have a mix of numeric data types (28 int64 and 11 float64) and 43 objects. Among the numeric variables, there are continuous (ex: Lot Frontage, Lot Area) and discrete (ex: OverallQual, OverallCond, Bmnt Full Bath). Nominal (ex: BmntQual, Foundation), ordinal (ex: PoolQC) and binary (ex: CentralAir) qualitative variables. There are also time variables such as Yr Sold. All the characteristics are detailed in the feature description section.

## Data Cleaning:

Regarding data cleaning, duplicates were searched for and deleted.

27 features had missing values with a number ranging from 1 to 2917. According to the description of the data, NA is used as qualifiers for many features to mean that there is no pool or basement for example. Or NA or NaN means a missing value for python. Therefore, for these features, we will replace NA with another qualifier such as No Pool for example to mean that the house does not have a swimming pool. Depending on the study of each variable, the missing values were replaced and some features have been removed. Indeed, Regarding 'Garage Cars' and 'Garage Area', it is possible that the data does not exist because the property does not have a garage or the property has a garage but the data is truly missing. Therefore we will proceed as follows: if there is no garage then 'Garage Cars' and 'Garage Area' equal to zero.

According to the graphical visualizations, some variables presented outliers (Gr Liv Area, Lot Area, Low Qual Fin SF). For some of them, the Z-score method was used. Most outliers have been removed. Others were checked and we reflected a really on the ground.

Some numerical features are actually really categorical features as MSSubClass and MoSold. We chose to transform them into nominal qualitative variables, as follows.

## Data Exploratory Analysis:

After cleaning the data, a visual exploration of the numerical and categorical data was carried out using pair plot and box plot. These visualizations allowed on the one hand, the validation of data cleaning. On the other hand, they made it possible to highlight polynomial relationships between certain features such as SalePrice and OverallQual or Year Built for example.

A correlation study was carried out on the numerical features. Disregarding the target, we see that the features are not too correlated with each other, it is a good thing. However, we note that 3 pairs of variables are strongly correlated with each other: 1st Flr SF/Total Bmnt SF, Gr Liv Area/TotRms AbvGr and Garage Area/Garage cars.

We will study the relationships between the qualitative features and the target using box plot, in order to carry out some hypothesis testing. The results of these tests showed significant differences between certain categories of characteristics and average house prices, such as MS Zoning and Central Air for example.

## Feature Engineering:

It was also studied the representation of the different modalities of the features. As a result, the 'Neighborhoods' variable has been modified for better representativeness of its modalities.

In some characteristics like 'Overall Qual' and 'Gr Liv Qual', we noticed an upward curved relationship rather than a simple linear relationship. Polynomial features were created for these characteristics.

In addition, interaction features have also been created for certain features such as 'Overall Qual' et 'Year Built' or 'Lot Area'.

Then, a One Hot encoding was performed on the categorical variables using Pandas one-hot encoder (get\_dummies).

## Modelisation:

The dataset was split into train set and test set, with a test size of 0.25%.

4 algorithms were tested on the dataset with a cross validation of 0.5%.

A simple linear regression was performed with GridSearchCV. Regularized linear regressions were also performed with RidgeCV, LassoCV and ElasticNetCV. For the latter 3, the numerical data was scaled using the StandardScaler transformer. The metrics studied during these modelisations were the Pearson coefficient (r2) and the root mean squared error (RMSE).

For these 4 regressions, very good results were obtained with r2 of 0.924 (simple linear regression) and 0.930 for the others. However, the lowest RMSE was obtained with LassoCV for a value of 22095. The best performing model on this dataset is LassoCV with 4-split cross-validation and standardization using StandardScaler.

## Improvement:

In order to improve the performance of these models:

- a further analysis of the functionalities could be considered, with hypothesis tests carried out on all the modalities of the categorical features for an optimized choice of the features, and with the establishment of other interactions of the features.
- The more in-depth study of the hyper-parameters could also contribute to the improvement of the model.

## A) Import the required libraries

```
In [1]: # Importing data
import pandas as pd
# Mathematical operations
import numpy as np
# For visualizing data
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.express as px
# For statistical analysis
from scipy import stats
from statsmodels.formula.api import ols
from statsmodels.stats.anova import anova_lm
# For splitting data
from sklearn.model_selection import train_test_split, KFold
# For scaling data
from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScaler
# For modelisation
from sklearn.linear_model import LinearRegression, LassoCV, RidgeCV, ElasticNetCV
from sklearn.model_selection import GridSearchCV
# Metrics: mse, rmse, r2, r2_adj, mean_squared_error, mean_absolute_error
```

## B) Dataset description and understanding data

The Ames housing dataset examines features of houses sold in Ames (a small city in the state of Iowa in the United States) during the 2006–2010 timeframe.

### Import dataset

```
In [2]: data = pd.read_csv("https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-ML0232EN-SkillsNetwork/ames_housing.txt")
Out[2]:
```

	Order	PID	SubClass	MS	Zoning	Frontage	Lot Area	Street	Alley	Lot Shape	Land Contour	...	Pool Area	Pool QC	Fence	Misc Feature	Misc Val	Mo Sold	Yr	
0	1	52630100	20	RL	1410	31770	Pave	NaN	NaN	IR1	Lvl	...	0	NaN	NaN	NaN	0	5	2010	
1	1	52630100	20	RL	1410	31770	Pave	NaN	NaN	IR1	Lvl	...	0	NaN	NaN	NaN	0	5	2010	
2	2	52635040	20	RH	810	14267	Pave	NaN	Reg	Lvl	...	...	0	NaN	MnPrv	NaN	0	6	2010	
3	3	52635103	20	RL	810	14267	Pave	NaN	IR1	Lvl	...	...	0	NaN	NaN	NaN	Gar2	12500	6	2010
4	4	526353030	20	RL	930	11160	Pave	NaN	Reg	Lvl	...	...	0	NaN	NaN	NaN	0	4	2010	

5 rows x 20 columns

### Feature description

```
In [4]: data_description = pd.read_csv('data_description.txt', sep=' ')
data_description
for k in data_description['MSSubClass'].identifies the type of dwelling involved in the sale.\t':
    print(k)

20      1-STORY 1946 & NEWER ALL STYLES
30      1-STORY 1945 & OLDER
40      1-STORY W/FINISHED ATTIC ALL AGES
45      1-1/2 STORY - UNFINISHED ALL AGES
50      2-STORY 1946 & NEWER ALL AGES
60      2-STORY 1946 & NEWER
70      2-STORY 1945 & OLDER
80      2-1/2 STORY ALL AGES
100     SPLIT OR MULTI-LEVEL
120     SPLIT FOYER
150     DUPLEX - ALL STYLES AND AGES
160     2-STORY 1946 & NEWER
180     PUD - NEW TIER LEVEL
190     2 FAMILY CONVERSION - ALL STYLES AND AGES
MSSZoning: Municipal zoning classification of the sale.
A      Agricultural
C      Commercial
RM      Residential Medium Density
RH      Residential High Density
RL      Residential Low Density
RP      Residential Low Density Park
RM      Residential Medium Density
LotFrontage: linear feet of street connected to property
LotArea: lot size in square feet
Street: Type of road access to property
Crty      Curved
Pave      Paved
Alley: Type of alley access to property
GrVl      Gravel
Pave      Paved
NA      No alley access
LotShape: General shape of property
Reg      Regular
IR1      Slightly irregular
IR2      Moderately irregular
IR3      Irregular
LandContour: Flatness of the property
Lvl      Level
HLS      Hillside - Significant slope from side to side
Low      Low
Utilities: Type of utilities available
AllPub      All public Utilities (E,G,W,& S)
NoSewer      Electricity, gas, & water (Septic Tank)
NoSewer      Electricity and Gas Only
ELO      Electricity only
LotConfig: Lot configuration
Inside      Inside lot
Corner      Corner lot
CulSacs      Cul-de-sacs
FR1      Frontage on 1 side of property
FR2      Frontage on 2 sides of property
FR3      Frontage on 3 sides of property
LandSlope: Slope of property
Gtl      Gentle slope
Mod      Moderate slope
Sev      Severe slope
Neighborhood: Physical locations within Ames city limits
Blmngtn      Bloomington Heights
Bluestem      Bluestem
BrDale      Briardale
BkSide      Brookside
ClearCr      Clear Creek
CollGr      College Creek
Crawfor      Crawford
Edwards      Edwards
Gilbert      Gilbert
IDOTRR      Iowa DOT and Rail Road
HwyEx      Highway Exchange
Mitchel      Mitchell
Names      Names
NorthAm      North Ames
NPkVill      Northpark Villa
Nridgdt      Northridge Heights
NWAmes      Northwest Ames
OldTwn      Old Town
SWISU      South & West of Iowa State University
Sawyer      Sawyer
SawyerW      Sawyer West
Somerset      Somerset
Stonerock      Stonerock
Timber      Timberland
Veenker      Veenker
Condition1: Proximity to various conditions
Atty      Adjacent to arterial street
Feedr      Adjacent to feeder street
Norm      Normal
RRNn      Within 200' of North-South Railroad
RRAn      Adjacent to North-South Railroad
PosA      Near positive off-site feature-park, greenbelt, etc.
PosN      Adjacent to positive off-site feature
RRMe      Within 200' of East-West Railroad
RRAn      Adjacent to East-West Railroad
Condition2: Proximity to various conditions (if more than one is present)
Artery      Adjacent to arterial street
Feedr      Adjacent to feeder street
Norm      Normal
RRNn      Within 200' of North-South Railroad
RRAn      Adjacent to North-South Railroad
PosA      Near positive off-site feature-park, greenbelt, etc.
PosN      Adjacent to positive off-site feature
RRMe      Within 200' of East-West Railroad
RRAn      Adjacent to East-West Railroad
BldgType: Type of dwelling
1Fam      Single-Family Detached
2FmCon      Two-family conversion; originally built as one-family dwelling
Duplx      Duplex
TwnHse      Townhouse End Unit
TwnHse1      Townhouse Inside Unit
HouseStyle: Style of dwelling
1Story      One story
1.5Fin      One and one-half story: 2nd level finished
1.5Unf      One and one-half story: 2nd level unfinished
2Story      Two story
2.5Fin      Two and one-half story: 2nd level finished
2.5Unf      Two and one-half story: 2nd level unfinished
SFoyer      Split Foyer
SLvl      Split Level
OverallQual: Rates the overall material and finish of the house
10      Very Excellent
9      Excellent
8      Very Good
7      Good
6      Above Average
5      Average
4      Below Average
3      Fair
2      Poor
1      Very Poor
OverallCond: Rates the overall condition of the house
10      Very Excellent
9      Excellent
8      Very Good
7      Good
6      Above Average
5      Average
4      Below Average
3      Fair
2      Poor
1      Very Poor
YearBuilt: Original construction date
YearRemodAdd: Remodel date (same as construction date if no remodeling or additions)
RoofStyle: Type of roof
Flat      Flat
Gable      Gable
Gambrel      Gambrel (Barn)
Hip      Hip
Mansard      Mansard
Shed      Shed
RoofMatl: Roof material
ClayTile      Clay or Tile
CompShg      Standard (Composite) Shingle
Hmbrn      Membrane
Metal      Metal
Roll      Roll
TarGrv      Gravel & Tar
WdShk     Wood Shakes
WdShnbl     Wood Shingles
Exterior1: Exterior covering on house
AsbShn     Asbestos Shingles
AsphShn     Asphalt Shingles
BrkComm     Brick Common
BrkFace     Brick Face
CBLOCK     Cinder Block
CmentBd     Cement Board
HdBdBoard     Hard Board
ImStucc     Imitation Stucco
MetalSd     Metal Siding
Other      Other
Plywood     Plywood
PreCast     Precast
Stone      Stone
Stucco      Stucco
VinylSd     Vinyl Siding
Wd Sdg      Wood Siding
WdShngl     Wood Shingles
Exterior2: Exterior covering on house (if more than one material)
AsbShn     Asbestos Shingles
AsphShn     Asphalt Shingles
BrkComm     Brick Common
BrkFace     Brick Face
CBLOCK     Cinder Block
CmentBd     Cement Board
HdBdBoard     Hard Board
ImStucc     Imitation Stucco
MetalSd     Metal Siding
Other      Other
Plywood     Plywood
PreCast     Precast
Stone      Stone
Stucco      Stucco
VinylSd     Vinyl Siding
Wd Sdg      Wood Siding
WdShngl     Wood Shingles
MasVnrType: Masonry veneer type
BrkCmn     Brick Common
BrkFace     Brick Face
CBLOCK     Cinder Block
None      None
Stone      Stone
MasFinArea: Masonry veneer area in square feet
ExterQual: Evaluates the quality of the material on the exterior
Ex      Excellent
Gd      Good
TA      Average/Typical
Fa      Fair
Po      Poor
ExterCond: Evaluates the present condition of the material on the exterior
Ex      Excellent
Gd      Good
TA      Average/Typical
Fa      Fair
Po      Poor
Foundation: Type of foundation
BkTtl     Brick & Tile
CBLOCK     Cinder Block
PCONE     Poured Concrete
Slab      Slab
2nd Flr     Second Floor
Stone      Stone
Wood      Wood
BmtQual: Evaluates the height of the basement
Ex      Excellent (>100' inches)
Gd      Good (80-99' inches)
TA      Typical (60-89' inches)
Fa      Fair (40-59' inches)
Po      Poor (<40' inches)
NA      No Basement
BmtCond: Evaluates the general condition of the basement
Ex      Excellent
Gd      Good
TA      Typical - slight dampness allowed
Fa      Fair - dampness or some cracking or settling
Po      Poor - Severe cracking, settling, or wetness
NA      No Basement
BmtExposure: Refers to walkout or garden level walls
Gd      Good Exposure
Av      Average Exposure (split levels or foyers typically score average or above)
Mn      Minimum Exposure
No      No Exposure
NA      No Basement
BmtFinType1: Rating of basement finished area
GLQ      Good Living Quarters
ALQ      Average Living Quarters
BLQ      Below Average Living Quarters
Rec      Average Rec Room
LwQ      Low Quality
Unf      Unfinished
NA      No Basement
BmtFinType2: Type 1 finished square feet
BmtFinType2: Rating of basement finished area (if multiple types)
GLQ      Good Living Quarters
ALQ      Average Living Quarters
BLQ      Below Average Living Quarters
Rec      Average Rec Room
LwQ      Low Quality
Unf      Unfinished
NA      No Basement
BmtUnfSF2: Type 2 finished square feet
BmtUnfSF2: Unfinished square feet of basement area
TotalBmtSF: Total square feet of basement area
Heating: Type of heating
Floor      Floor
GasV      Gas forced warm air furnace
GASW      Gas hot water or steam heat
Grav      Gravity
OthW      Hot water or steam heat other than gas
Wall      Wall
HeatingQC: Heating quality and condition
Ex      Excellent
Gd      Good
TA      Average/Typical
Fa      Fair
Po      Poor
CentralAir: Central air conditioning
N      No
Y      Yes
Electrical: Electrical system
SBrkr      Standard Circuit Breakers & Romex
FuseA      Fuse Box over 60 AMP and all Romex wiring (Average)
FuseF      60 AMP Fuse Box and mostly Romex wiring (Fair)
FuseP      60 AMP Fuse Box and mostly knob & tube wiring (poor)
MiscVal: Miscellaneous feature not covered in other categories
1stFlrSF: First floor square feet
2ndFlrSF: Second floor square feet
LowQualFinSF: Low quality finished square feet (all floors)
GrLivArea: Above grade (ground) living area square feet
BmtFullBath: Basement full bathrooms
BmtHalfBath: Basement half bathrooms
FullBath: Full bathrooms above grade
HalfBath: Half baths above grade
Bedroom: bedroom above grade (does NOT include basement bedrooms)
Kitchen: Kitchens above grade
KitchenQual: Kitchen quality
Ex      Excellent
Gd      Good
TA      Typical/Average
Fa      Fair
Po      Poor
TotRmsAbvGrd: Total rooms above grade (does not include bathrooms)
Functional: Home functionality (Assume typical unless deductions are warranted)
Typ      Typical
Min1     Minor Deductions 1
Min2     Minor Deductions 2
Mod      Moderate Deductions
Major1   Major Deductions 1
Major2   Major Deductions 2
Sev      Severely Damaged
Salvage  Salvage only
Fireplaces: Number of fireplaces
FireplaceQu: Fireplace quality
Ex      Excellent - Exceptional Masonry Fireplace
Gd      Good - Masonry Fireplace in main level
TA      Average - Prefabricated Fireplace in main living area or Masonry Fireplace in basement
Fa      Fair - Prefabricated Fireplace in basement
Po      Poor - Ben Franklin Stove
NA      No Fireplace
GarageType: Garage location
2Type      More than one type of garage
Attchd     Attached to home
Basement     Basement (Garage part of house - typically has room above garage)
BuiltIn     Built-in garage
CarPort     Car Port
Detached     Detached from home
NA      No Garage
GarageYrBlt: Year garage was built
GarageFinish: Interior finish of the garage
Fin      Finished
Rfrn      Rough finished
Unf      Unfinished
NA      No Garage
GarageCars: Size of garage in car capacity
GarageArea: Size of garage in square feet
GarageQual: Garage quality
Ex      Excellent
Gd      Good
TA      Typical/Average
Fa      Fair
Po      Poor
NA      No Garage
GarageCond: Garage condition
Ex      Excellent
Gd      Good
TA      Typical/Average
Fa      Fair
Po      Poor
NA      No Garage
PavedDrive: Paved driveway
Y      Yes
P      Partial Paveement
N      No
WoodDeckSF: Wood deck area in square feet
OpenPorchSF: Open porch area in square feet
EnclosedPorch: Enclosed porch area in square feet
3SeasonPorch: Three season porch area in square feet
ScreenPorch: Screen porch area in square feet
PoolArea: Pool area in square feet
PoolQC: Pool quality
Ex      Excellent
Gd      Good
TA      Average/Typical
Fa      Fair
Po      Poor
NA      No Pool
Fence: Fence quality
GdPrv     Good Privacy
MnPrv     Minimum Privacy
GdWo      Good Wood/Wire
MnWo      Minimum Wood/Wire
NA      No Fence
MiscFeatures: Miscellaneous feature not covered in other categories
Elv      Elevator
Gar2      2nd Garage (if not described in garage section)
Othr      Other
Shed      Shed (over 100 SF)
TennisC     Tennis Court
NA      None
MiscVal: Value of miscellaneous feature
MSold: Month Sold (MM)
YrSold: Year Sold (YYYY)
SaleType: Type of sale
WD      Warranty Deed - Conventional
CWD      Warranty Deed - Cash
Vndr      Vendor
New      Home just constructed and sold
COD      Court Order/Deed/Broker
Con      Contract 15% Down payment regular terms
ConLw     Contract Low Down payment and low interest
ConLI     Contract Low Interest
ConLD     Contract Low Down
Oth      Other
SaleCondition: Condition of sale
Normal      Normal Sale
Abnorml     Abnormal Sale - trade, foreclosure, short sale
AdjLand     Adjoining land purchase
Alloca      Allocation - two linked properties with separate deeds, typically condo with a garage unit
Partly      Partly between family members
Partial     Home was not completed when last assessed (associated with New Homes)
```

## Dataset information

```
In [5]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2931 entries, 0 to 2930
Data columns (total 82 columns):
#   Column                Non-Null Count  Dtype
---  --
0   Order                 2931 non-null   int64
1   PID                   2931 non-null   int64
2   MS SubClass           2931 non-null   int64
3   MS Zoning              2931 non-null   object
4   Lot Frontage          2441 non-null   float64
5   Lot Area              2931 non-null   int64
6   Street                2931 non-null   object
7   Alley                198 non-null    object
8   Lot Shape             2931 non-null   object
9   Land Contour          2931 non-null   object
10  Utilities             2931 non-null   object
14  Condition 2          2931 non-null   object
16  Bldg Type             2931 non-null   object
17  House Style           2931 non-null   object
18  Overall Qual          2931 non-null   int64
19  Overall Cond          2931 non-null   int64
20  Year Built            2931 non-null   int64
21  Year Remod/Add        2931 non-null   int64
22  Roof Style            2931 non-null   object
23  Roof Matl             2931 non-null   object
24  Exterior 1st          2931 non-null   object
25  Exterior 2nd          2931 non-null   object
26  Mas Vnr Type          2908 non-null   object
27  Mas Vnr Area          2908 non-null   float64
28  Exter Qual            2931 non-null   object
29  Exter Cond            2931 non-null   object
30  Foundation            2931 non-null   object
31  Bmnt Qual             2851 non-null   object
32  Bmnt Exposure         2848 non-null   object
33  Bmnt Fin Type 1       2851 non-null   object
34  Bmnt Fin SF 1         2930 non-null   float64
35  Bmnt Fin Type 2       2850 non-null   object
36  Bmnt Fin SF 2         2930 non-null   float64
37  Bmnt Unf SF           2930 non-null   float64
38  Total Bmnt SF         2930 non-null   float64
40  Heating              2931 non-null   object
41  Heating QC            2931 non-null   object
42  Central Air           2930 non-null   object
43  Electrical            2930 non-null   object
44  1st Flr SF            2931 non-null   int64
45  2nd Flr SF            2931 non-null   int64
46  Low Qual Fin SF       2931 non-null   int64
47  Gr Liv Area           2931 non-null   int64
48  Bmnt Full Bath        2829 non-null   float64
49  Bmnt Half Bath        2829 non-null   float64
50  Full Bath             2931 non-null   int64
51  Half Bath             2931 non-null   int64
52  Bedroom AbvGrd       2931 non-null   int64
53  Kitchen Qual          2931 non-null   object
54  Kitchen Area          2931 non-null   int64
55  TotRms AbvGrd        2931 non-null   int64
56  Functional            2931 non-null   object
57  Fireplaces            2931 non-null   int64
58  Fireplace Qu          1509 non-null   int64
59  Garage Type           2774 non-null   object
60  Garage Yr Blt         2772 non-null   float64
61  Garage Finish         2932 non-null   object
62  Garage Cars           2930 non-null   float64
63  Garage Area           2930 non-null   float64
64  Garage Qual           2932 non-null   object
65  Garage Cond           2772 non-null   object
66  Paved Drive           2931 non-null   object
67  Wood Deck SF         2931 non-null   int64
68  Open Porch SF        2931 non-null   int64
69  Enclosed Porch        2931 non-null   int64
70  3Season Porch         2931 non-null   int64
71  Screen Porch          2931 non-null   int64
72  Pool Area             2931 non-null   int64
73  Pool QC              13 non-null    object
74  Fence                 172 non-null   object
75  Misc Feature          1166 non-null   object
76  Misc Val              2931 non-null   int64
77  Mo Sold               2931 non-null   int64
78  Yr Sold               2931 non-null   int64
79  Sale Type             2931 non-null   object
80  Sale Condition        2931 non-null   object
81  Sale Price            2931 non-null   int64
dtypes: float64(11), int64(28), object(43)
memory usage: 1.8+ MB
```

- According to the output above, we have 2931 entries, from 0 to 2930, along with 81 features.
- We notice that there are missing values for some features.
- 'SalePrice' is our target or response variable and the rest of the characteristics are our predictor variables.
- We also have a mix of numeric data types (28 int64 and 11 float64) and 43 objects.

Among the numeric variables, there are continuous (ex: Lot Frontage, Lot Area) and discrete (ex: OverallQual, OverallCond, Bmnt Full Bath, Nominal (ex: BmntQual, Foundation), ordinal (ex: PoolQC) and binary (ex: CentralAir) qualitative variables. There are also time variables such as Yr Sold.

## Analysis of response variable

```
In [6]: data['SalePrice'].describe()
Out[6]:
```

	count	2931	0.000000
mean	160807.729167		
std	79875.557267		
min	12789.000000		
25%	129500.000000		
50%	160000.000000		
75%	213500.000000		
max	755000.000000		
Name:	SalePrice,	dtype: float64	

From the above analysis, it is important to note that the minimum value is greater than 0. Also, there is a large difference between the minimum value and the 25th percentile. It is greater than the 75th percentile and the maximum value. This means that our data may not be normally distributed (an important assumption for linear regression analysis), so will check for normality in the Log Transform section.

## Looking for Correlations

Before proceeding with data cleaning, it is useful to establish a correlation between the target variable (SalePrice) and other predictor variables, since some of them might not have a major impact in determining the price of the house and not to be used in the analysis.

There are many ways to discover the correlation between the target variable and the rest of the characteristics. Pair plots, scatter plots, heat maps and correlation matrices are the most common. Below, we will use the corr() function to list the top features based on Pearson's correlation coefficients. We also have a few closely two number sequences are <0.1 (too low). Correlation coefficient can only be calculated on numeric attributes (float and integer), therefore only numeric attributes will be selected.

We will select correlations greater than 0.5.

```
In [7]: data_num = data.select_dtypes(include = ('float64', 'int64'))
data_num.columns = data_num.columns[['SalePrice']]
best_num_corr = data_num.corr(abs(data_num.corr)>0.5).sort_values(ascending=False)
print(f'The {len(best_num_corr)} most highly correlated variables with SalePrice:')
best_num_corr
Out[7]:
```

	Overall Qual	0.795225
1	Gr Liv Area	0.647891
2	Garage Cars	0.647891
3	Garage Area	0.647891
4	Total Bmnt SF	0.632775
5	1st Flr SF	0.621672
6	Year Built	0.553340
7	Garage Yr Blt	0.532809
8	Mass Vnr Area	0.509277
Name:	SalePrice,	dtype: float64

Next, let's generate a few per plot to visually inspect the correlation between some of these features and the target variable. We will use the seaborn sns.pairplot() function for this analysis. Also, constructing pair plots is one of the possible ways to spot outliers that might be present in the data.

```
In [8]: for i in range(1, len(data_num.columns), 5):
sns.pairplot(data[data_num.columns[i:i+5], y_vars = ['SalePrice']])
```

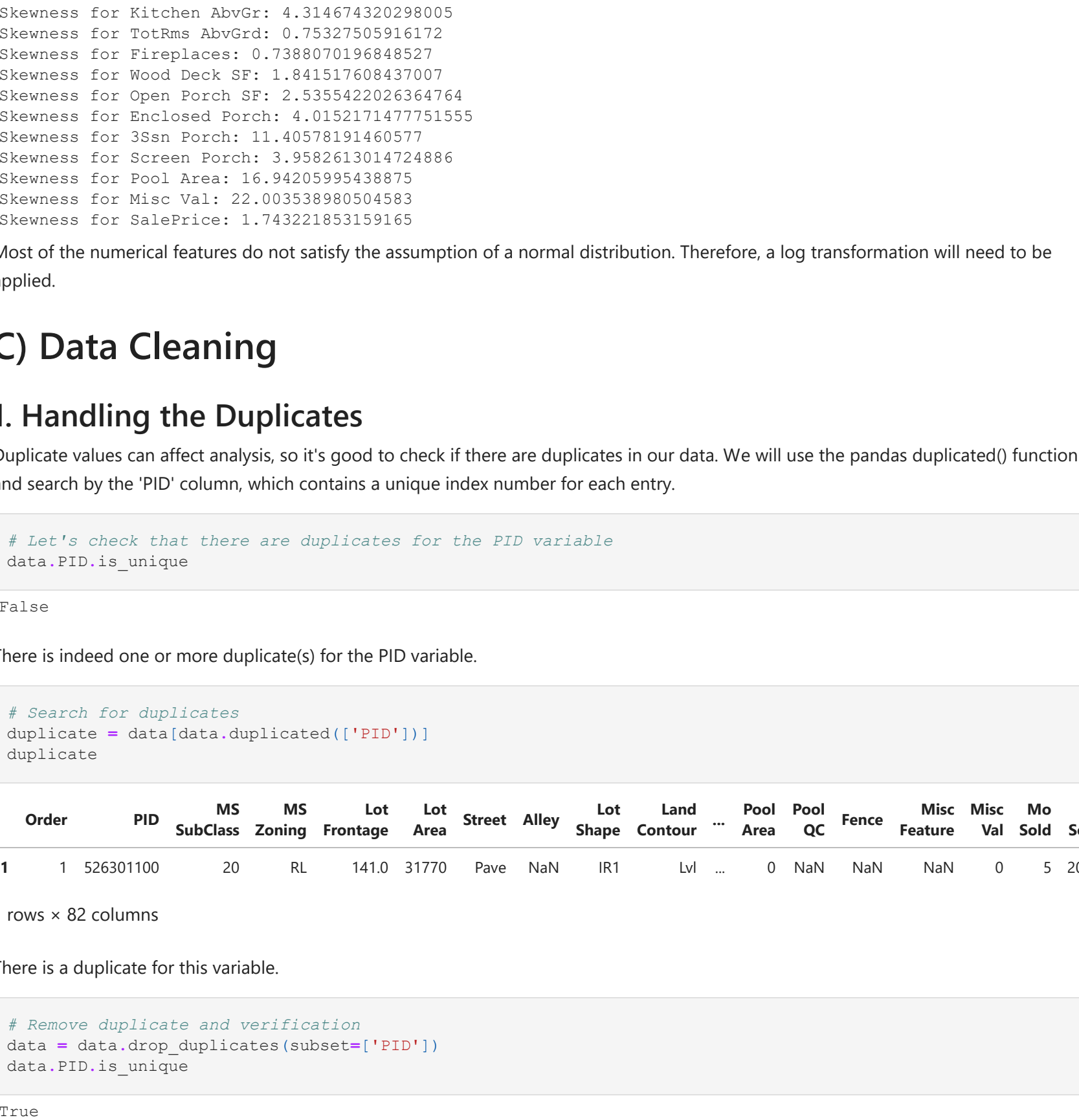




Visually, we notice that Overall Qual, Gr Liv Area, Garage Cars, Garage Area and more features are correlated to the target 'SalePrice'. We also see that there are possibly outliers for some features like Lot Area, 1st Flr SF and Gr Liv Area.

The heat maps method also makes it possible to check whether certain features are not correlated with each other, which could bias the model.

```
In [9]: data_num_corr = data_num_corr()
heatmap_num_corr = data_num_corr[abs(data_num_corr)>=0.5]
plt.figure(figsize=(15, 10))
sns.heatmap(heatmap_num_corr,
            annot=True,
            fmt='.1f',
            cmap='magma',
            vmin=-0.5,
            vmax=1)
title = plt.title('Correlations between numerical features')
```



It can be seen that some features are strongly correlated with each other, such as Garage Yr Bld and Year Built or Gr Liv Area and TotRms AbvGrd for example. A features selection will certainly be required for modeling.

## Feature distribution

The normal distribution assumption must be satisfied in order to perform any type of regression analysis. There are several ways to verify this assumption. Here, we can simply use the skew() function to calculate our level of feature skewness. We consider that the distribution is not normally distributed if the level of asymmetry is less than 0.5 or greater than 0.5.

```
In [10]: for i in data_num_columns:
skewness = data[i].skew()
if skewness < -0.5 or skewness > 0.5:
print(f'Skewness for {i}: {skewness}')
```

Skewness for MS SubClass: 1.3579444255012858  
Skewness for Lot Frontage: 1.49877456787124  
Skewness for Lot Area: 12.778040640965061  
Skewness for Overall Cond: 0.5749245539249781  
Skewness for Year Built: -0.6041550817882403  
Skewness for Mas Vnr Area: 2.6073702204244156  
Skewness for BsmFin SF 1: 1.4118746240939502  
Skewness for BsmFin SF 2: 4.140793742804169  
Skewness for BsmUnf SF: 0.9234449946334007  
Skewness for Total Bsm SF: 1.1536330901822214  
Skewness for 1st Flr SF: 1.467869250607155  
Skewness for 2nd Flr SF: 0.8669705520648831  
Skewness for Low Qual Fin SF: 12.12027034881609  
Skewness for Gr Liv Area: 1.27395787171539  
Skewness for Bsm Full Bath: 0.615697638287629  
Skewness for Bsm Half Bath: 3.94192438497531  
Skewness for Half Bath: 0.6992547390083218  
Skewness for Kitchen AbvGrd: 4.314674320298005  
Skewness for TotRms AbvGrd: 0.7520730591172  
Skewness for Fireplaces: 0.7388070196848527  
Skewness for Wood Deck SF: 1.841517608437007  
Skewness for Open Porch SF: 3.9582613014724866  
Skewness for Enclosed Porch: 4.015217147771555  
Skewness for 3Ssn Porch: 11.40578191460577  
Skewness for Screen Porch: 3.9582613014724866  
Skewness for Pool Area: 16.94205995438975  
Skewness for Misc Val: 22.005358980504583  
Skewness for SalePrice: 1.743021853153165

Most of the numerical features do not satisfy the assumption of a normal distribution. Therefore, a log transformation will need to be applied.

## C) Data Cleaning

### 1. Handling the Duplicates

Duplicate values can affect analysis, so it's good to check if there are duplicates in our data. We will use the pandas duplicated() function and search by the 'PID' column, which contains a unique index number for each entry.

```
In [11]: # Let's check that there are duplicates for the PID variable
data.PID.is_unique

Out[11]: False
```

There is indeed one or more duplicate(s) for the PID variable.

```
In [12]: # Search for duplicates
duplicate = data[data.duplicated(['PID'])]
duplicate

Out[12]:
```

Order	PID	MS SubClass	Zoning	Frontage	Lot Area	Street	Alley	Lot Shape	Land Contour	Utilities	Lot Config	Pool Area	Pool QC	Fence	Misc Feature	Misc Val	Mo Sold	Yr Sold	Yr Sold	Yr Sold
1	1	526301100	20	RL	141.0	31770	Pave	NaN	IR1	Lvl	...	0	NaN	NaN	NaN	0	5	2010		

1 rows × 22 columns

There is a duplicate for this variable.

```
In [13]: # Remove duplicate and verification
data = data.drop_duplicates(subset=['PID'])
data.PID.is_unique

Out[13]: True
```

### 2. Handling the Missing Values

#### Missing values research

Missing data can significantly bias modeling. We create a function to determine the percentage of missing values according to the different predictive variables.

```
In [14]: def missing_data(data):
missing_data_count = data.isna().sum()
missing_data_percent = data.isna().sum() / len(data) * 100

missing_data = pd.DataFrame({
    'Number_Na': missing_data_count,
    '%_NaNa': missing_data_percent
})

missing_data = missing_data[missing_data['Number_Na'] > 0]
missing_data.sort_values(by='Number_Na', ascending=False, inplace=True)

print(f'There are {missing_data.shape[0]} features with missing values.\n')
return missing_data
```

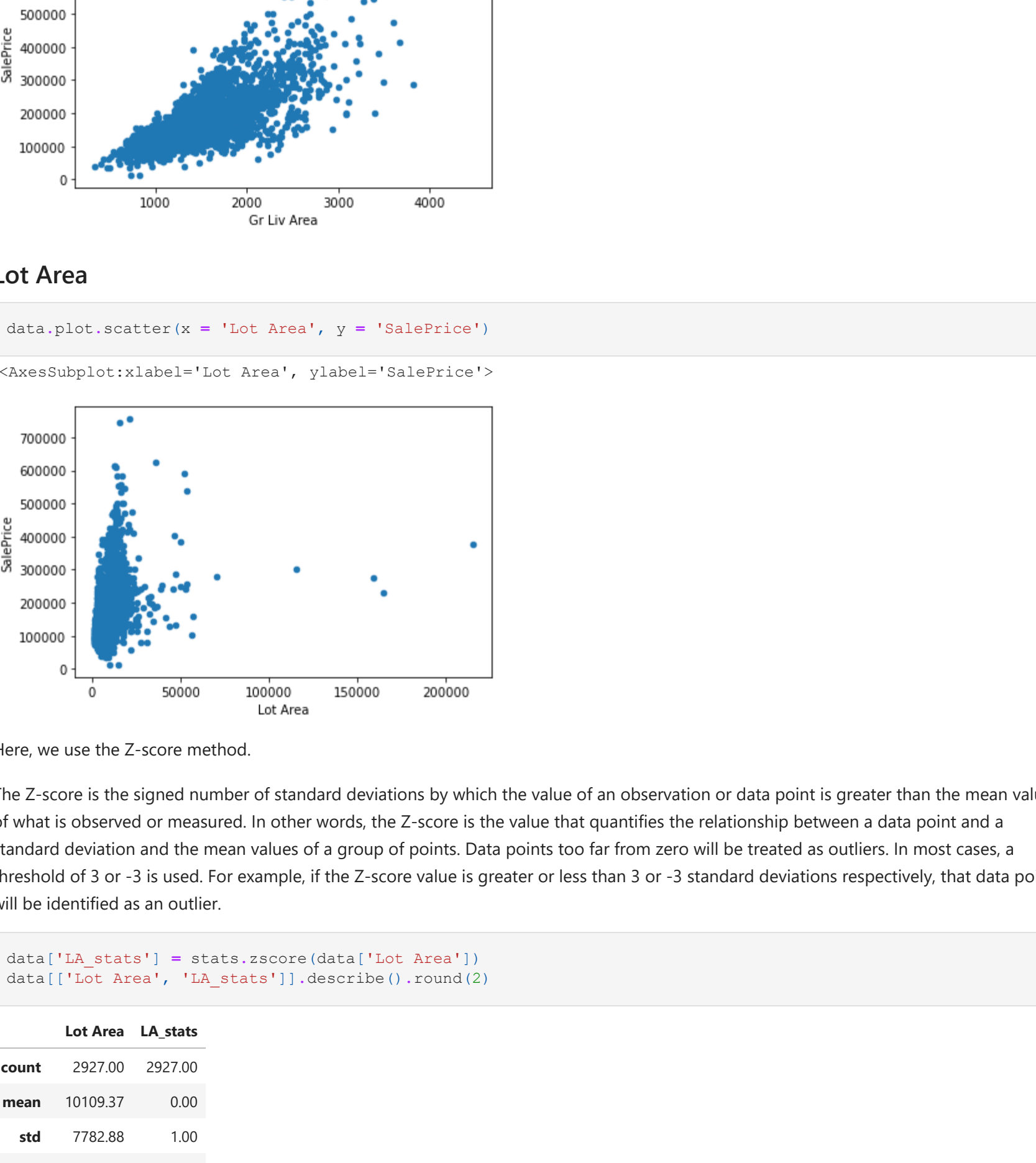
missing\_data(data)

There are 27 features with missing values.

	Number_Na	%_NaNa
Pool QC	2917	99.556314
Misc Feature	2824	96.382253
Alley	2732	93.242321
Fence	2358	80.477816
Fireplace Qu	1422	48.532423
Lot Frontage	490	16.723549
Garage Cond	159	5.426621
Garage Qual	159	5.426621
Garage Finish	159	5.426621
Garage Yr Bld	159	5.426621
Garage Type	157	5.358362
Bsmnt Exposure	83	2.832765
Bsmnt Fin Type 2	81	2.764505
Bsmnt Fin Type 1	80	2.730375
Bsmnt Qual	80	2.730375
Bsmnt Type	80	2.730375
Mas Vnr Area	23	0.784983
Mas Vnr Type	23	0.784983
Bsmnt Half Bath	2	0.068259
Bsmnt Full Bath	2	0.068259
Total Bsmnt SF	1	0.034130
Bsmnt Unf SF	1	0.034130
Garage Cars	1	0.034130
Garage Area	1	0.034130
Bsmnt Fin SF 2	1	0.034130
Bsmnt Fin SF 1	1	0.034130
Electrical	1	0.034130

```
In [16]: missing_values = data.isnull().sum().sort_values(ascending=False)
missing_values = missing_values.head(27)
missing_values.plot(kind='bar', figsize=(15,6), fontsize=10)
plt.xlabel('Features', fontsize=20)
plt.ylabel('Number', fontsize=20)
plt.title('Missing values', fontsize=25)
```

Text(0.5, 1.0, 'Missing values')



### Replacement of missing values and deletion of variables

According to the description of the data, NA is used as qualifiers for many features to mean that there is no pool or basement for example. Gr Liv Area or NaH means a missing value for python. Therefore, for these features, we will replace the NA with another qualifier such as No Pool for example to mean that the house does not have a swimming pool.

```
In [17]: data['Pool QC'] = data['Pool QC'].replace(np.nan, 'No Pool')
data['Misc Feature'] = data['Misc Feature'].replace(np.nan, 'None')
data['Alley'] = data['Alley'].replace(np.nan, 'No Alley')
data['Fence'] = data['Fence'].replace(np.nan, 'No Fence')
data['Fireplace Qu'] = data['Fireplace Qu'].replace(np.nan, 'No Fireplace')
data['Garage Cond'] = data['Garage Cond'].replace(np.nan, 'No Garage')
data['Garage Finish'] = data['Garage Finish'].replace(np.nan, 'No Garage')
data['Garage Qual'] = data['Garage Qual'].replace(np.nan, 'No Garage')
data['Garage Type'] = data['Garage Type'].replace(np.nan, 'No Garage')
data['Bsmnt Exposure'] = data['Bsmnt Exposure'].replace(np.nan, 'No Basement')
data['Bsmnt Fin Type 2'] = data['Bsmnt Fin Type 2'].replace(np.nan, 'No Basement')
data['Bsmnt Fin Type 1'] = data['Bsmnt Fin Type 1'].replace(np.nan, 'No Basement')
data['Bsmnt Qual'] = data['Bsmnt Qual'].replace(np.nan, 'No Basement')
data['Bsmnt Cond'] = data['Bsmnt Cond'].replace(np.nan, 'No Basement')
data['Bsmnt Full Bath'] = data['Bsmnt Full Bath'].replace(np.nan, 'No Basement')
data['Bsmnt Unf SF'] = data['Bsmnt Unf SF'].replace(np.nan, 'No Basement')
data['Mas Vnr Type'] = data['Mas Vnr Type'].replace(np.nan, 'None')
```

missing\_data(data)

There are 12 features with missing values.

	Number_Na	%_NaNa
Lot Frontage	490	16.723549
Garage Yr Bld	159	5.426621
Mas Vnr Area	23	0.784983
Bsmnt Full Bath	2	0.068259
Bsmnt Half Bath	2	0.068259
Bsmnt Fin SF 1	1	0.034130
Bsmnt Fin SF 2	1	0.034130
Bsmnt Unf SF	1	0.034130
Total Bsmnt SF	1	0.034130
Electrical	1	0.034130
Garage Cars	1	0.034130
Garage Area	1	0.034130

Regarding 'Garage Cars' and 'Garage Area', it is possible that the data does not exist because the property does not have a garage or the property has a garage but the data is truly missing. Therefore we will proceed as follows: if there is no garage then 'Garage Cars' and 'Garage Area' equal to zero.

```
In [19]: data['Garage Cars'] = np.where(data['Garage Cond'] == 'No Garage', 0, data['Garage Cars'])
data['Garage Area'] = np.where(data['Garage Cond'] == 'No Garage', 0, data['Garage Area'])
```

Regarding 'Garage Yr Bld', we choose to delete the entire column, because this variable is strongly correlated to x and cannot be replaced by any date.

```
In [20]: data = data.drop('Garage Yr Bld', axis=1)
```

Similarly, if the house has no basement, we will set the following variables to zero.

```
In [21]: data['Bsmnt Full Bath'].unique().tolist()
```

[1.0, 0.0, 2.0, 2.0, 3.0, nan]

```
In [22]: data['Bsmnt Fin SF 1'] = np.where(data['Bsmnt Qual'] == 'No Basement', 0, data['Bsmnt Fin SF 1'])
data['Bsmnt Fin SF 2'] = np.where(data['Bsmnt Qual'] == 'No Basement', 0, data['Bsmnt Fin SF 2'])
data['Bsmnt Unf SF'] = np.where(data['Bsmnt Qual'] == 'No Basement', 0, data['Bsmnt Unf SF'])
data['Total Bsmnt SF'] = np.where(data['Bsmnt Qual'] == 'No Basement', 0, data['Total Bsmnt SF'])
data['Bsmnt Full Bath'] = np.where(data['Bsmnt Qual'] == 'No Basement', 0, data['Bsmnt Full Bath'])
data['Bsmnt Half Bath'] = np.where(data['Bsmnt Qual'] == 'No Basement', 0, data['Bsmnt Half Bath'])
```

# Verification

```
data['Bsmnt Full Bath'].unique().tolist()
```

[1.0, 0.0, 2.0, 3.0]

Similarly, if the property does not have a type of masonry veneer, the Masonry veneer area is considered equal zero square feet.

```
In [24]: data['Mas Vnr Area'] = np.where(data['Mas Vnr Type'] == 'None', 0, data['Mas Vnr Area'])
```

For electricity, we choose to replace the missing value with the most frequent value in the data, which is Standard Circuit Breakers & Romex.

```
In [25]: data['Electrical'].value_counts()

Out[25]:
Electrical    2682
FuseA         188
FuseB          50
FuseF           8
Misc           1
Name: Electrical, dtype: int64
```

```
In [26]: data['Electrical'].fillna('SBKr', inplace = True)
```

Regarding the Linear feet of street connected to property, we replace missing values by 0.

```
In [27]: data['Lot Frontage'] = data['Lot Frontage'].fillna(0)
```

### Verification of missing values

```
In [28]: missing_data(data)

Out[28]:
```

There are 0 features with missing values.

	Number_Na	%_NaNa
Lot Frontage	490	16.723549
Garage Yr Bld	159	5.426621
Mas Vnr Area	23	0.784983
Bsmnt Full Bath	2	0.068259
Bsmnt Half Bath	2	0.068259
Bsmnt Fin SF 1	1	0.034130
Bsmnt Fin SF 2	1	0.034130
Bsmnt Unf SF	1	0.034130
Total Bsmnt SF	1	0.034130
Electrical	1	0.034130
Garage Cars	1	0.034130
Garage Area	1	0.034130

We can remove the Order and PID columns which will bring nothing to our predictions.

```
In [29]: data.drop(['PID', 'Order'], axis=1, inplace=True)
```

## 3. Handling the Outliers

According to the graphical analysis carried out in 'looking correlation' part, there are some outliers. Outliers can significantly affect our models and can be a valuable source of information, providing us with insights into specific behaviors.

### Gr Liv Area

```
In [30]: data.plot.scatter(x = 'Gr Liv Area', y = 'SalePrice')
```

<AxesSubplot: xlabel='Gr Liv Area', ylabel='SalePrice'>



There are 3 values that deviate from the trend (Gr Liv Area > 4000 and SalePrice < 300000). We will delete them.

```
In [31]: GrLivArea_outlier = data[(data['Gr Liv Area'] > 4000) & (data['SalePrice'] < 300000)]
GrLivArea_outlier
```

	MS SubClass	MS Zoning	Lot Area	Lot Area	Street	Alley	Lot Shape	Land Contour	Utilities	Lot Config	Pool Area	Pool QC	Fence	Misc Feature	Misc Val	Mo Sold	Yr Sold	Yr Sold	Yr Sold
1499	60	RL	3130	63887	Pave	No Alley	IR3	Low	AllPub	Corner	...	480	Gd	No Fence	None	0	6	2009	WD
2181	20	RL	1280	39290	Pave	No Alley	IR1	Bnk	AllPub	Inside	...	0	No Pool	No Fence	Elev	17000	10	2007	
2182	60	RL	1300	40094	Pave	No Alley	IR1	Bnk	AllPub	Inside	...	0	No Pool	No Fence	None	0	10	2007	

3 rows × 20 columns

```
In [32]: data = data.drop(GrLivArea_outlier.index)
```

```
In [33]: # Verification
data.plot.scatter(x = 'Gr Liv Area', y = 'SalePrice')
```



### Lot Area

```
In [34]: data.plot.scatter(x = 'Lot Area', y = 'SalePrice')
```



Here, we use the Z-score method.

The Z-score is the signed number of standard deviations by which the value of an observation or data point is greater than the mean value of what is observed or measured. In other words, the Z-score is the value that quantifies the relationship between a data point and a standard deviation and the mean values of a group of points. Data points too far from zero will be treated as outliers. In most cases, a threshold of 3 or -3 is used. For example, if the Z-score value is greater or less than 3 or -3 standard deviations respectively, that data point will be identified as an outlier.

```
In [35]: data['LA_stats'] = stats.zscore(data['Lot Area'])
data[['Lot Area', 'LA_stats']].describe().round(2)
```

	Lot Area	LA_stats
count	2927.00	2927.00
mean	10109.37	0.00
std	7782.88	1.00
min	1300.00	-1.13
25%	7439.00	-0.34
50%	9430.00	-0.09
75%	11523.00	0.18
max	215245.00	26.36

The scaled results show a mean of 0.000 and a standard deviation of 1.000, indicating that the transformed values fit the z-scale model. The maximum value of 26.36 is further evidence of the presence of outliers, as it is well above the z-score limit of +3.

We choose to delete the 4 points whose Lot Area is greater than 100000.

```
In [36]: LA_outlier = data.sort_values(by = 'Lot Area', ascending=False)[4:]
LA_outlier
```

	MS SubClass	MS Zoning	Lot Area	Lot Area	Street	Alley	Lot Shape	Land Contour	Utilities	Lot Config	Pool Area	Pool QC	Fence	Misc Feature	Misc Val	Mo Sold	Yr Sold	Yr Sold	Yr Sold
957	20	RL	1502	215245	Pave	No Alley	IR3	Low	AllPub	Inside	...	0	No Pool	No Fence	None	0	6	2009	WD
1571	190	RL	0.0	164660	Grvl	No Alley	IR1	HLS	AllPub	Corner	...	0	No Pool	No Fence	Shed	700	8	2008	WD
2116	50	RL	0.0	159000	Pave	No Alley	IR2	Low	AllPub	CudDsc	...	0	No Pool	No Fence	Shed	500	6	2007	WD
2072	20	RL	0.0	115149	Pave	No Alley	IR2	Low	AllPub	CudDsc	...	0	No Pool	No Fence	None	0	6	2007	WD

4 rows × 20 columns

```
In [37]: data = data.drop(LA_outlier.index)
```

```
In [38]: # Verification
data.plot.scatter(x = 'Lot Area', y = 'SalePrice')
```



In [39]: data['LA\_stats'] = stats.zscore(data['Lot Area'])
data[['Lot Area', 'LA\_stats']].describe().round(2)

	Lot Area	LA_stats
count	2923.00	2923.00
mean	9899.45	0.00
std	5154.64	1.00
min	1300.00	-1.67
25%	7437.00	-0.48
50%	9416.00	-0.09
75%	11506.00	0.31
max	70761.00	11.79

According to the Z-score method, there are still outliers but we consider that these remaining values reflect a reality on the ground.

### Low Qual Fin SF

```
In [40]: data.plot.scatter(x = 'Low Qual Fin SF', y = 'SalePrice')
```



In [41]: data[['LQFSF\_State', 'LA\_stats']] = stats.zscore(data[['Low Qual Fin SF']])
data[['Low Qual Fin SF', 'LQFSF\_State']].describe().round(2)

	Low Qual Fin SF	LQFSF_State
count	2923.00	2923.00
mean	4.69	0.00
std	46.37	1.00
min	0.00	-0.10
25%	0.00	-0.10
50%	0.00	-0.10
75%	0.00	-0.10
max	1064.00	22.85

The result shows the presence of outliers. Here, we choose to delete the last value, namely the largest.

```
In [42]: LQFSF_outlier = data.sort_values(by = 'Low Qual Fin SF', ascending=False)[1:]
LQFSF_outlier
```

	MS SubClass	MS Zoning	Lot Area	Lot Area	Street	Alley	Lot Shape	Land Contour	Utilities	Lot Config	Pool Area	Pool QC	Fence	Misc Feature	Misc Val	Mo Sold	Yr Sold	Yr Sold	Yr Sold
661	50	RL	60.0	10410	Pave	Grvl	Reg	Lvl	AllPub	Inside	...	0	No Pool	No Fence	None	0	6	2009	WD

1 rows × 20 columns

```
In [43]: data = data.drop(LQFSF_outlier.index)
```

```
In [44]: # Verification
data.plot.scatter(x = 'Low Qual Fin SF', y = 'SalePrice')
```



We delete the 2 variables that were used to study the Z-score.

```
In [45]: data.drop(['LQFSF_State', 'LA_stats'], axis=1, inplace=True)
```

## 4. Feature modification







As is, each feature is treated as an independent quantity. However, there may be **interaction effects**, in which the impact of one feature may depend on the current value of a different feature.

For example, there may be a higher premium for increasing 'Overall Qual' for houses that were built more recently. If such a premium or a similar effect exists, a feature that multiplies 'Overall Qual' by 'Year Built' can help us capture it.

Another style of interaction term involves feature proportions: for example, to get at something like quality per square foot we could divide 'Overall Qual' by 'Lot Area'.

```
In [76]: data['OQ_X_YB'] = data['Overall Qual'] * data['Year Built']
data['OQ_/_LA'] = data['Overall Qual'] / data['Lot Area']
```

### 3. Feature Scaling

A transformer such as `StandardScaler()`, `MinMaxScaler()` or `RobustScaler()` will be used later on the train and test set in order to scale the numerical features.

### 4. One-hot encoding for dummy variables

A significant challenge, especially when dealing with data with many columns, is ensuring that each column is correctly encoded. This is especially true with columns of data that are ordered categories (ordinals) versus unordered categories. Unordered categories should be one-hot encoded, but this can greatly increase the number of features and create features that are highly correlated with each other, which is not necessarily the case here.

The Pandas one-hot encoder (`get_dummies`) works well with data defined as categorical. By default it only converts object columns. We will also delete the first column.

However, one-hot encoding can overfit data with much higher mean squared errors on test data than no encoding.

```
In [77]: cat_cols = data.select_dtypes(include = ['object']).columns

print(f'There are {len(cat_cols.tolist())} object variables to encode')

cat_cols.tolist()

There are 45 object variables to encode

Out[77]: ['MS_SubClass',
'MS_Zoning',
'Street',
'Alley',
'Lot_Shape',
'Land_Contour',
'Utilities',
'Lot_Config',
'Land_Slope',
'Neighborhood',
'Condition 1',
'Condition 2',
'Bidg Type',
'House Style',
'Roof Style',
'Roof Matl',
'Exterior 1st',
'Exterior 2nd',
'Mas Vnr Type',
'Exter Qual',
'Exter Cond',
'Foundation',
'Bsmt Qual',
'Bsmt Cond',
'Bsmt Exposure',
'BsmtFin Type 1',
'BsmtFin Type 2',
'Heating',
'Heating QC',
'Central Air',
'Electrical',
'Functional',
'Fireplace Qu',
'Garage Type',
'Garage Finish',
'Garage Qual',
'Garage Cond',
'Paved Drive',
'Pool QC',
'Fence',
'Misc Feature',
'Mo Sold',
'Sale Type',
'Sale Condition']

In [78]: data[cat_cols].apply(lambda x: x.nunique()).sort_values(ascending=False)

Neighborhood      26
Exterior 2nd      17
MS_SubClass       16
Exterior 1st      16
Mo Sold           12
Sale Type         10
Condition 1        9
Functional          8
Condition 2        8
House Style        8
Garage Type        7
BsmtFin Type 2     7
Roof Matl          7
BsmtFin Type 1     7
MS_Zoning          7
Fireplace Qu       6
Foundation         6
Heating            6
Garage Qual        6
Garage Cond        6
Bsmt Cond          6
Bsmt Qual          6
Sale Condition     6
Roof Style         6
Electrical         5
Misc Feature       5
Fence              5
Pool_QC            5
Kitchen Qual       5
Exter Cond         5
Heating QC         5
Bidg Type          5
Mas Vnr Type       5
Bsmt Exposure      5
Lot_Config         5
Exter Qual         4
Garage Finish      4
Land Contour       4
Lot_Shape          4
Land_Slope         3
Utilities           3
Paved Drive        3
Alley              3
Central Air        2
Street             2
Style             t64

In [79]: data[cat_cols].apply(lambda x: x.nunique()).sort_values(ascending=False).sum()

306
```

No categorical variable has only one value. We can therefore encode them all. The 45 object columns will become 306 columns in total. We'll use the Pandas One Hot Encoder (`get_dummies`), which by default only converts object columns.

```
In [80]: #data.info()

In [81]: data = pd.get_dummies(data)

In [82]: data.shape

Out[82]: (2922, 344)
```

### Summary EDA

After all the processing provided, the dataset now has 2922 rows and 344 columns.

- Duplicates have been removed.
- Outliers were removed.
- Missing values have been processed.
- The distribution of variables has been studied and corrected.
- Hypothesis tests were carried out.
- New features such as polynomials have been created and others have been modified.
- All the categorical features were encoded.

## G) Modelisation

### 1. Train test split

```
In [83]: print('data shape:', data.shape)

X = data.drop('SalePrice', axis = 1)
y = data['SalePrice']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 5)

print('X_train shape:', X_train.shape)
print('X_test shape:', X_test.shape)
print('y_train shape:', y_train.shape)

data.shape      (2922, 344)
X_train.shape    (2191, 343)
X_test.shape     (731, 343)
y_train.shape    (2191,)
y_test.shape     (731,)
```

### 2. Simple linear regression with GridSearchCV

```
In [84]: param_grid = {'fit_intercept': ['bool', 'True'],
                    'n_jobs': np.arange(1,21)}

grid = GridSearchCV(LinearRegression(), param_grid, cv=4, scoring='r2')

grid.fit(X_train, y_train)

print(grid.best_score_)
print(grid.best_params_)
print(grid.best_estimator_)

model_lr = grid.best_estimator_

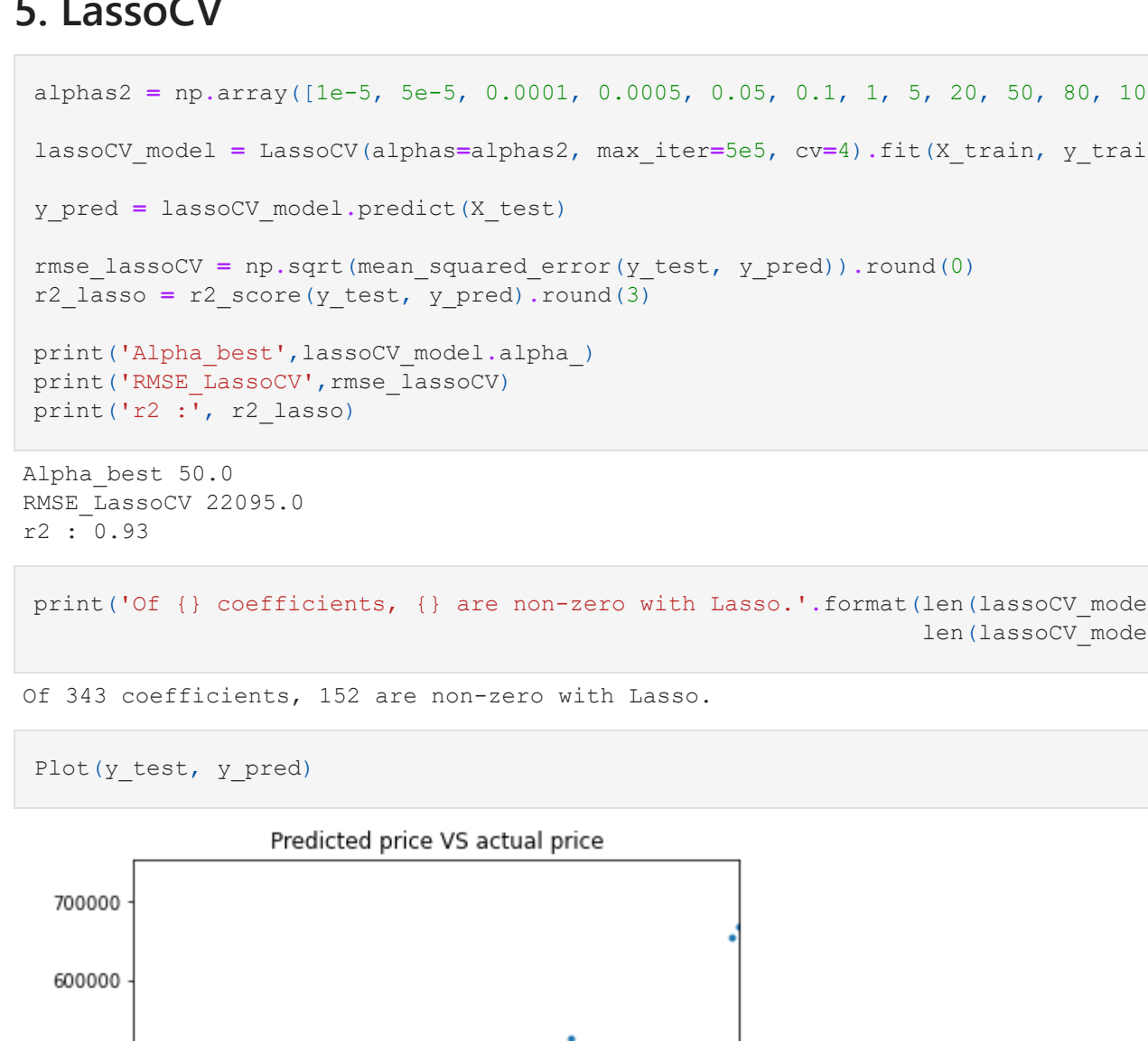
y_pred = model_lr.predict(X_test)

r2_lr = model_lr.score(X_test, y_test).round(3)
rmse_lr = np.sqrt(mean_squared_error(y_test, y_pred)).round(0)
print('RMSE_LinearRegression:',rmse_lr)
print('r2 :', r2_lr)

0.9229727482155242
{'fit_intercept': 'bool', 'n_jobs': 1}
LinearRegression(fit_intercept='bool', n_jobs=1)
RMSE_LinearRegression: 23050.0
r2 : 0.924
```

```
In [85]: def Plot(y_test, y_pred):
    f = plt.figure(figsize=(6,6))
    ax = plt.axes()
    ax.plot(y_test, y_pred, marker='o', ls='', ms=3.0)
    ax.annotate("r-squared = {:.3f}".format(r2_score(y_test, y_pred)), (0,1))
    lim = (0, y_test.max())
    ax.set(xlabel='Actual Price',
           ylabel='Predicted Price',
           xlim=lim,
           ylim=lim,
           title='Predicted price VS actual price', );

In [86]: Plot(y_test, y_pred)
```



### 3. Scaling numerical features

```
In [87]: scaler = StandardScaler()

cat_col = X_train.select_dtypes(include = ['float64', 'int64']).columns
X_train[num_col] = scaler.fit_transform(X_train[num_col])
X_test[num_col] = scaler.transform(X_test[num_col])
```

### 4. RidgeCV

```
In [88]: alphas = [0.005, 0.05, 0.1, 0.3, 1, 3, 5, 10, 15, 30, 80]

ridgecv_model = RidgeCV(alphas=alphas, cv=4).fit(X_train, y_train)

y_pred = ridgecv_model.predict(X_test)

rmse_ridgecv = np.sqrt(mean_squared_error(y_test, y_pred)).round(0)
r2_ridgecv = r2_score(y_test, y_pred).round(3)

print('Alpha_best',ridgecv_model.alpha_)
print('RMSE_RidgeCV',rmse_ridgecv)
print('r2 :', r2_ridgecv)

Alpha_best 10.0
RMSE_RidgeCV 22149.0
r2 : 0.93

In [89]: Plot(y_test, y_pred)
```



### 5. LassoCV

```
In [90]: alphas2 = np.array([1e-5, 5e-5, 0.0001, 0.0005, 0.05, 0.1, 1, 5, 20, 50, 80, 100, 120, 140])

lassocv_model = LassoCV(alphas=alphas2, max_iter=5e5, cv=4).fit(X_train, y_train)

y_pred = lassoCV_model.predict(X_test)

rmse_lassocv = np.sqrt(mean_squared_error(y_test, y_pred)).round(0)
r2_lassocv = r2_score(y_test, y_pred).round(3)

print('Alpha_best',lassocv_model.alpha_)
print('RMSE_LassoCV',rmse_lassocv)
print('r2 :', r2_lassocv)

Alpha_best 50.0
RMSE_LassoCV 22095.0
r2 : 0.93

In [91]: print('Of {} coefficients, {} are non-zero with Lasso.'.format(len(lassocv_model.coef_),
                                                                    len(lassocv_model.coef_[0]))

Of 343 coefficients, 152 are non-zero with Lasso.

In [92]: Plot(y_test, y_pred)
```



### 6. ElasticNetCV

```
In [93]: l1_ratios = np.linspace(0.1, 0.9, 9)

elasticnetcv_model = ElasticNetCV(alphas=alphas2,
                                  l1_ratio=l1_ratios,
                                  max_iter=15e5).fit(X_train, y_train)

y_pred = elasticnetcv_model.predict(X_test)

rmse_elasticnetcv = np.sqrt(mean_squared_error(y_test, y_pred)).round(0)
r2_elasticnetcv = r2_score(y_test, y_pred).round(3)

print('Alpha_best',elasticnetcv_model.alpha_)
print('RMSE_elasticnetcv',rmse_elasticnetcv)
print('r2 :', r2_elasticnetcv)

C:\Users\godet\anaconda3\lib\site-packages\sklearn\linear_model\_coordinate_descent.py:633: ConvergenceWarning:
Objective did not converge. You might want to increase the number of iterations. Duality gap: 8033072842.539082
model = cd_fast.enet_coordinate_descent_gra
C:\Users\godet\anaconda3\lib\site-packages\sklearn\linear_model\_coordinate_descent.py:633: ConvergenceWarning:
Objective did not converge. You might want to increase the number of iterations. Duality gap: 299155944.07861
3, tolerance: 1079843679.0488164
model = cd_fast.enet_coordinate_descent_gra
C:\Users\godet\anaconda3\lib\site-packages\sklearn\linear_model\_coordinate_descent.py:633: ConvergenceWarning:
Objective did not converge. You might want to increase the number of iterations. Duality gap: 2627979628.62359
6, tolerance: 1118639322.64597
model = cd_fast.enet_coordinate_descent_gra
C:\Users\godet\anaconda3\lib\site-packages\sklearn\linear_model\_coordinate_descent.py:633: ConvergenceWarning:
Objective did not converge. You might want to increase the number of iterations. Duality gap: 1424847463.729980
7, tolerance: 1047154139.775393
model = cd_fast.enet_coordinate_descent_gra
Alpha_best 0.05
RMSE_elasticnetcv 22150.0
r2 : 0.93
```

### 7. Result and discussion

```
In [94]: rmse_vals = [rmse_lr, rmse_ridgecv, rmse_lassocv, rmse_elasticnetcv]
r2_vals = [r2_lr, r2_ridgecv, r2_lassocv, r2_elasticnetcv]
labels = ['Linear', 'Ridge', 'Lasso', 'ElasticNet']

df = pd.Series(rmse_vals, index=labels).to_frame()
df.rename(columns={0: 'RMSE'}, inplace=1)
df['R2'] = pd.Series(r2_vals, index=labels).to_frame()
df
```

```
Out[94]:
```

	RMSE	R2
Linear	23050.0	0.924
Ridge	22149.0	0.930
Lasso	22095.0	0.930
ElasticNet	22150.0	0.930

```
In [95]: f = plt.figure(figsize=(10,10))
ax = plt.axes()

labels = ['Ridge', 'Lasso', 'ElasticNet']

models = [RidgeCV_model, lassoCV_model, elasticnetcv_model]

for mod, lab in zip(models, labels):
    ax.plot(y_test, mod.predict(X_test),
            marker='o', ls='', ms=3.0, label=lab)

X = data.drop('SalePrice', axis = 1)
y = data['SalePrice']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 5)
ax.plot(y_test, y_train, y_test, y_test, y_test, y_test,
        marker='o', ls='', ms=3.0, label='Linear')

leg = plt.legend(frameon=False)
leg.get_frame().set_edgecolor('black')
leg.get_frame().set_linewidth(1.0)

ax.set(xlabel='Actual Price',
       ylabel='Predicted Price',
       title='Predicted vs Actual Price')
```



For these 4 regressions, very good results were obtained with r2 of 0.924 (simple linear regression) and 0.93 for the others. However, the lowest RMSE was obtained with LassoCV for a value of 22095. The best performing model on this dataset is LassoCV with 4 splits cross-validation and standardization using StandardScaler.

```
In [ ]:
```