

In [1]:

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Comment this if the data visualisations doesn't work on your side
%matplotlib inline

plt.style.use('bmh')
```

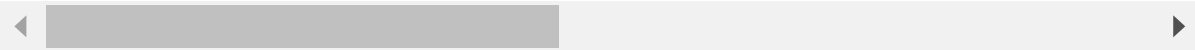
In [2]:

```
df = pd.read_csv("train.csv")
df.head()
```

Out[2]:

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

5 rows × 11 columns



In [3]:

df.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

```

In [4]:

```

# df.count() does not include NaN values
df2 = df[[column for column in df if df[column].count() / len(df) >= 0.3]]
del df2['Id']
print("List of dropped columns:", end=" ")
for c in df.columns:
    if c not in df2.columns:
        print(c, end=", ")
print('\n')
df = df2

```

List of dropped columns: Id, Alley, PoolQC, Fence, MiscFeature,

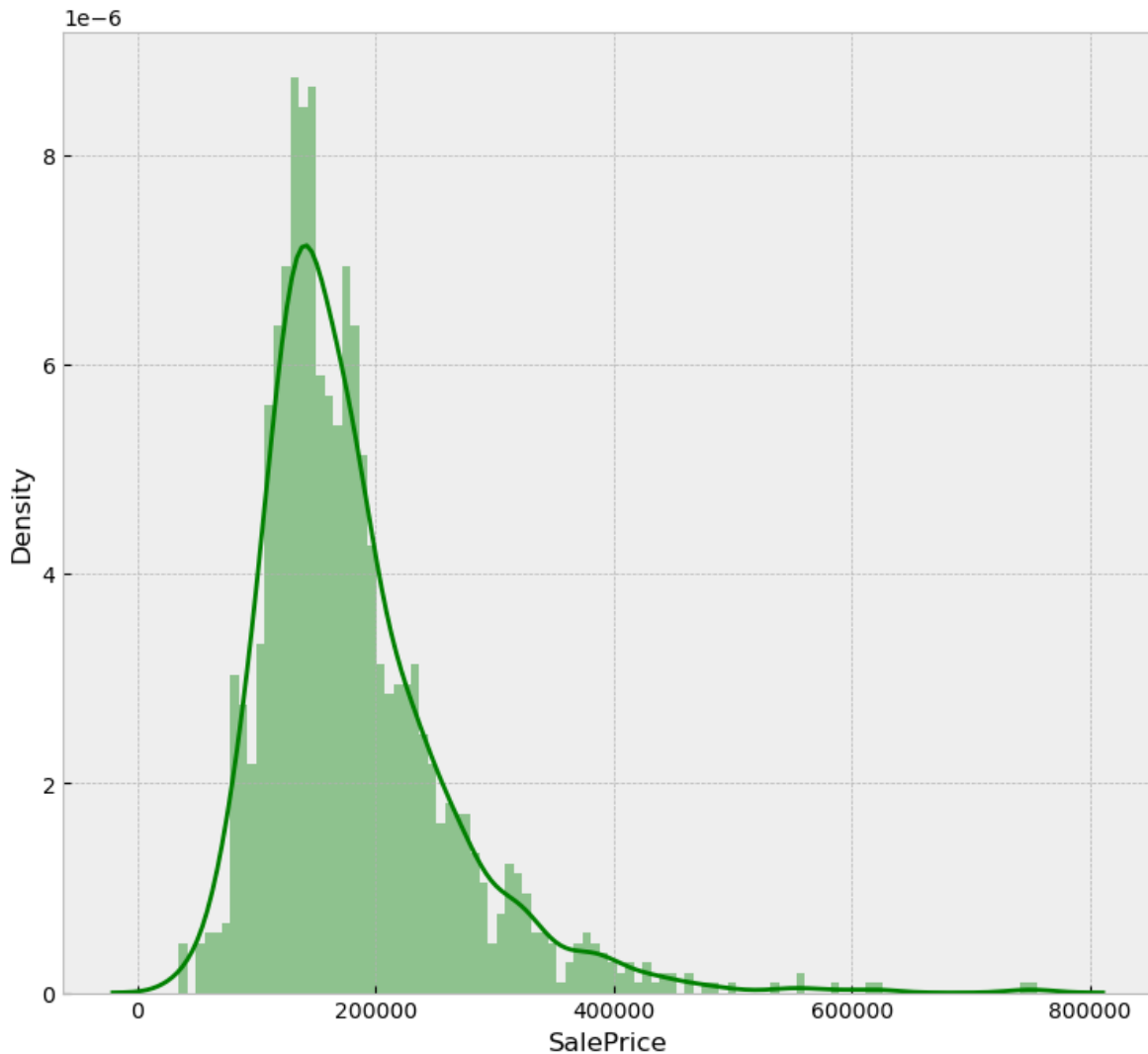
In [5]:

```
print(df['SalePrice'].describe())
plt.figure(figsize=(9, 8))
sns.distplot(df['SalePrice'], color='g', bins=100, hist_kws={'alpha': 0.4});
```

```
count      1460.000000
mean       180921.195890
std        79442.502883
min        34900.000000
25%       129975.000000
50%       163000.000000
75%       214000.000000
max        755000.000000
Name: SalePrice, dtype: float64
```

C:\Users\hp\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```



In [8]:

```
list(set(df.dtypes.tolist()))
```

Out[8]:

```
[dtype('O'), dtype('float64'), dtype('int64')]
```

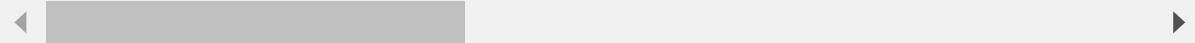
In [9]:

```
df_num = df.select_dtypes(include = ['float64', 'int64'])  
df_num.head()
```

Out[9]:

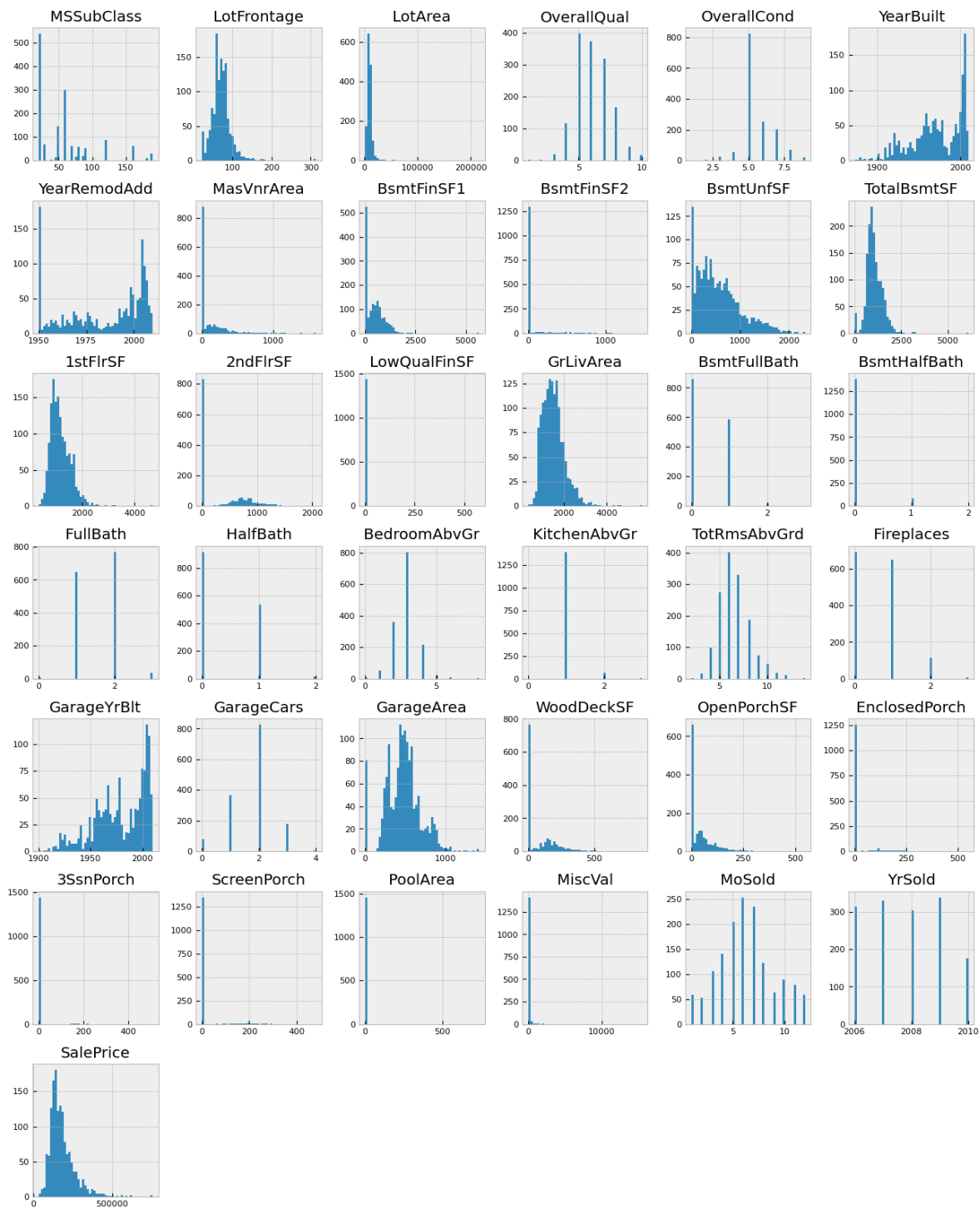
	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	Ma
0	60	65.0	8450	7	5	2003	2003	
1	20	80.0	9600	6	8	1976	1976	
2	60	68.0	11250	7	5	2001	2002	
3	70	60.0	9550	7	5	1915	1970	
4	60	84.0	14260	8	5	2000	2000	

5 rows × 37 columns



In [10]:

```
df_num.hist(figsize=(16, 20), bins=50, xlabelsize=8, ylabelsize=8);
```



In [11]:

```
df_num_corr = df_num.corr()['SalePrice'][:-1] # -1 because the latest row is SalePrice
golden_features_list = df_num_corr[abs(df_num_corr) > 0.5].sort_values(ascending=False)
print("There is {} strongly correlated values with SalePrice:\n{}".format(len(golden_features_list), golden_features_list))
```

There is 10 strongly correlated values with SalePrice:

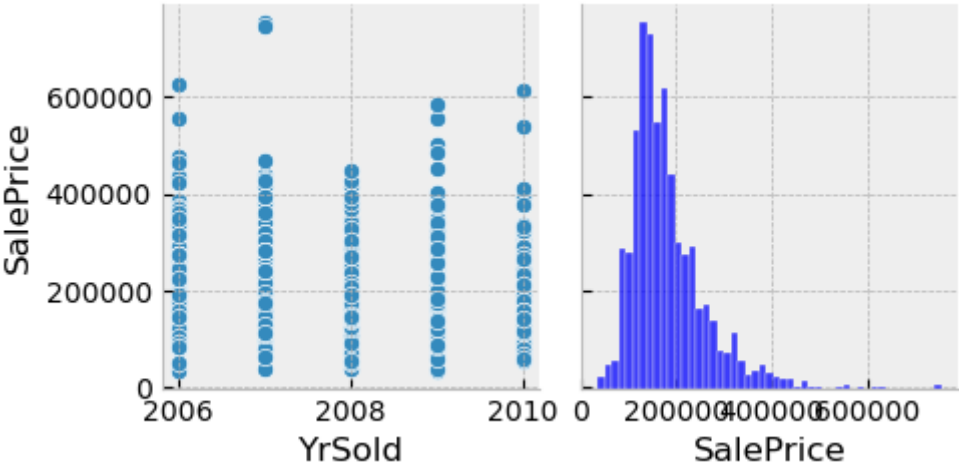
OverallQual	0.790982
GrLivArea	0.708624
GarageCars	0.640409
GarageArea	0.623431
TotalBsmntSF	0.613581
1stFlrSF	0.605852
FullBath	0.560664
TotRmsAbvGrd	0.533723
YearBuilt	0.522897
YearRemodAdd	0.507101

Name: SalePrice, dtype: float64

In [12]:

```
for i in range(0, len(df_num.columns), 5):
    sns.pairplot(data=df_num,
                  x_vars=df_num.columns[i:i+5],
                  y_vars=['SalePrice'])
```





In [13]:

```

import operator

individual_features_df = []
for i in range(0, len(df_num.columns) - 1): # -1 because the last column is SalePrice
    tmpDf = df_num[[df_num.columns[i], 'SalePrice']]
    tmpDf = tmpDf[tmpDf[df_num.columns[i]] != 0]
    individual_features_df.append(tmpDf)

all_correlations = {feature.columns[0]: feature.corr()['SalePrice'][0] for feature in indiv
all_correlations = sorted(all_correlations.items(), key=operator.itemgetter(1))
for (key, value) in all_correlations:
    print("{:>15}: {:>15}".format(key, value))

```

```

KitchenAbvGr: -0.1392006921778576
  HalfBath: -0.08439171127179902
  MSSubClass: -0.08428413512659509
  OverallCond: -0.07785589404867797
    YrSold: -0.028922585168736813
BsmtHalfBath: -0.02883456718548182
  PoolArea: -0.014091521506356765
BsmtFullBath: 0.011439163340408606
  MoSold: 0.046432245223819446
  3SsnPorch: 0.06393243256889088
  OpenPorchSF: 0.08645298857147718
  MiscVal: 0.08896338917298921
  Fireplaces: 0.12166058421363891
  BsmtUnfSF: 0.16926100049514173
  BedroomAbvGr: 0.18093669310848806
  WoodDeckSF: 0.1937060123752066
  BsmtFinSF2: 0.19895609430836594
EnclosedPorch: 0.24127883630117497
  ScreenPorch: 0.2554300795487841
    LotArea: 0.2638433538714051
  LowQualFinSF: 0.30007501655501323
  LotFrontage: 0.35179909657067737
  MasVnrArea: 0.43409021975689227
  BsmtFinSF1: 0.47169042652357296
  GarageYrBlt: 0.4863616774878596
  YearRemodAdd: 0.5071009671113866
    YearBuilt: 0.5228973328794967
  TotRmsAbvGrd: 0.5337231555820284
    FullBath: 0.5745626737760822
    1stFlrSF: 0.6058521846919153
  GarageArea: 0.6084052829168346
  TotalBsmtSF: 0.6096808188074374
  GarageCars: 0.6370954062078923
    2ndFlrSF: 0.6733048324568376
  GrLivArea: 0.7086244776126515
  OverallQual: 0.7909816005838053

```

In [14]:

```
golden_features_list = [key for key, value in all_correlations if abs(value) >= 0.5]
print("There is {} strongly correlated values with SalePrice:\n{}".format(len(golden_features_list), golden_features_list))
```

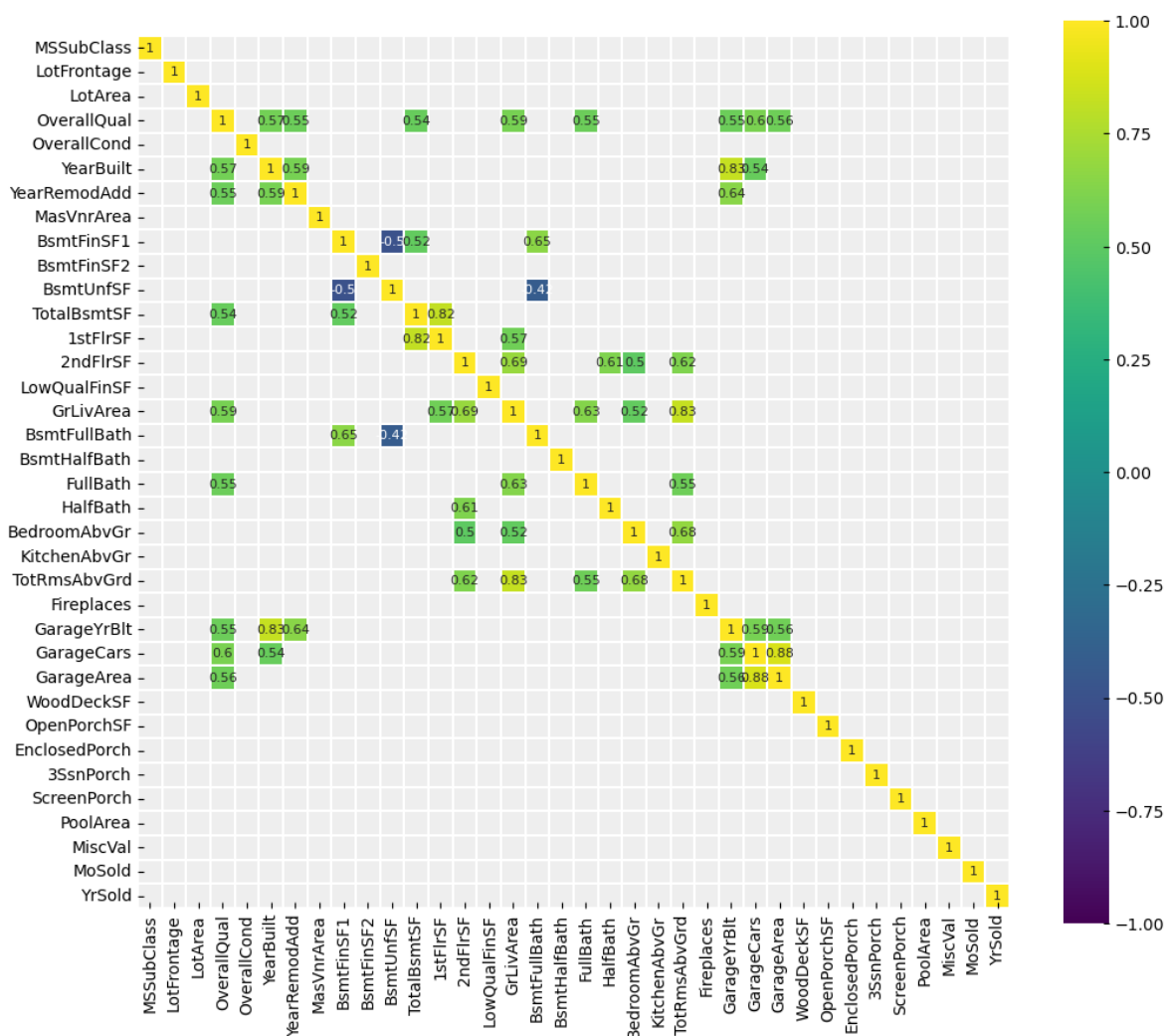
There is 11 strongly correlated values with SalePrice:

```
['YearRemodAdd', 'YearBuilt', 'TotRmsAbvGrd', 'FullBath', '1stFlrSF', 'GarageArea', 'TotalBsmtSF', 'GarageCars', '2ndFlrSF', 'GrLivArea', 'OverallQual']
```

In [15]:

```
corr = df_num.drop('SalePrice', axis=1).corr() # We already examined SalePrice correlations
plt.figure(figsize=(12, 10))
```

```
sns.heatmap(corr[(corr >= 0.5) | (corr <= -0.4)],
             cmap='viridis', vmax=1.0, vmin=-1.0, linewidths=0.1,
             annot=True, annot_kws={"size": 8}, square=True);
```



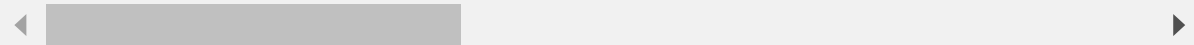
In [16]:

```
quantitative_features_list = ['LotFrontage', 'LotArea', 'MasVnrArea', 'BsmtFinSF1', 'BsmtFi
    '2ndFlrSF', 'LowQualFinSF', 'GrLivArea', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath', 'H
    'BedroomAbvGr', 'KitchenAbvGr', 'TotRmsAbvGrd', 'Fireplaces', 'GarageCars', 'GarageArea
    'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'PoolArea', 'MiscVal', 'SalePrice']
df_quantitative_values = df[quantitative_features_list]
df_quantitative_values.head()
```

Out[16]:

	LotFrontage	LotArea	MasVnrArea	BsmtFinSF1	BsmtFinSF2	TotalBsmtSF	1stFlrSF	2ndFlr
0	65.0	8450	196.0	706	0	856	856	8
1	80.0	9600	0.0	978	0	1262	1262	
2	68.0	11250	162.0	486	0	920	920	8
3	60.0	9550	0.0	216	0	756	961	7
4	84.0	14260	350.0	655	0	1145	1145	10

5 rows × 28 columns



In [17]:

```
features_to_analyse = [x for x in quantitative_features_list if x in golden_features_list]
features_to_analyse.append('SalePrice')
features_to_analyse
```

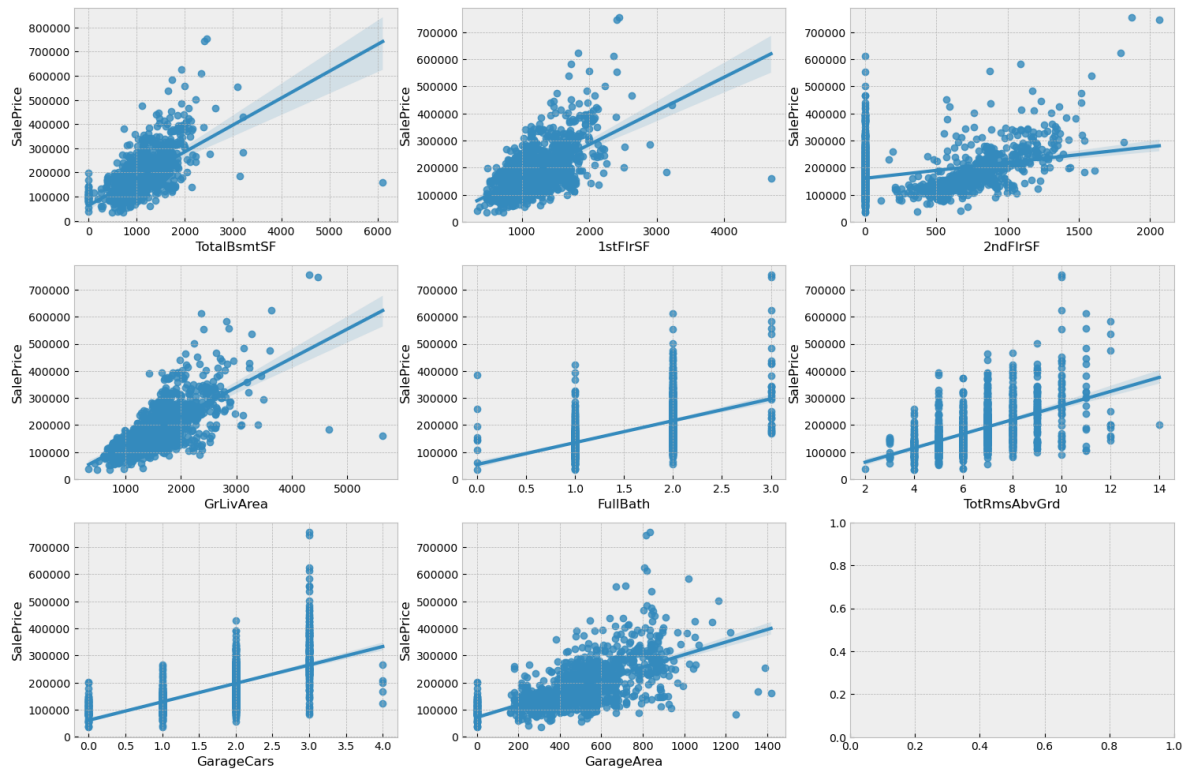
Out[17]:

```
['TotalBsmtSF',
 '1stFlrSF',
 '2ndFlrSF',
 'GrLivArea',
 'FullBath',
 'TotRmsAbvGrd',
 'GarageCars',
 'GarageArea',
 'SalePrice']
```

In [18]:

```
fig, ax = plt.subplots(round(len(features_to_analyse) / 3), 3, figsize = (18, 12))

for i, ax in enumerate(fig.axes):
    if i < len(features_to_analyse) - 1:
        sns.regplot(x=features_to_analyse[i],y='SalePrice', data=df[features_to_analyse], a
```



In [19]:

```
# quantitative_features_list[:-1] as the last column is SalePrice and we want to keep it
categorical_features = [a for a in quantitative_features_list[:-1] + df.columns.tolist() if
df_categ = df[categorical_features]
df_categ.head()
```

Out[19]:

	MSSubClass	MSZoning	Street	LotShape	LandContour	Utilities	LotConfig	LandSlope	Nei
0	60	RL	Pave	Reg	Lvl	AllPub	Inside	Gtl	
1	20	RL	Pave	Reg	Lvl	AllPub	FR2	Gtl	
2	60	RL	Pave	IR1	Lvl	AllPub	Inside	Gtl	
3	70	RL	Pave	IR1	Lvl	AllPub	Corner	Gtl	
4	60	RL	Pave	IR1	Lvl	AllPub	FR2	Gtl	

5 rows × 49 columns

In [20]:

```
df_not_num = df_categ.select_dtypes(include = ['O'])
print('There is {} non numerical features including:\n{}'.format(len(df_not_num.columns), d
```

There is 39 non numerical features including:

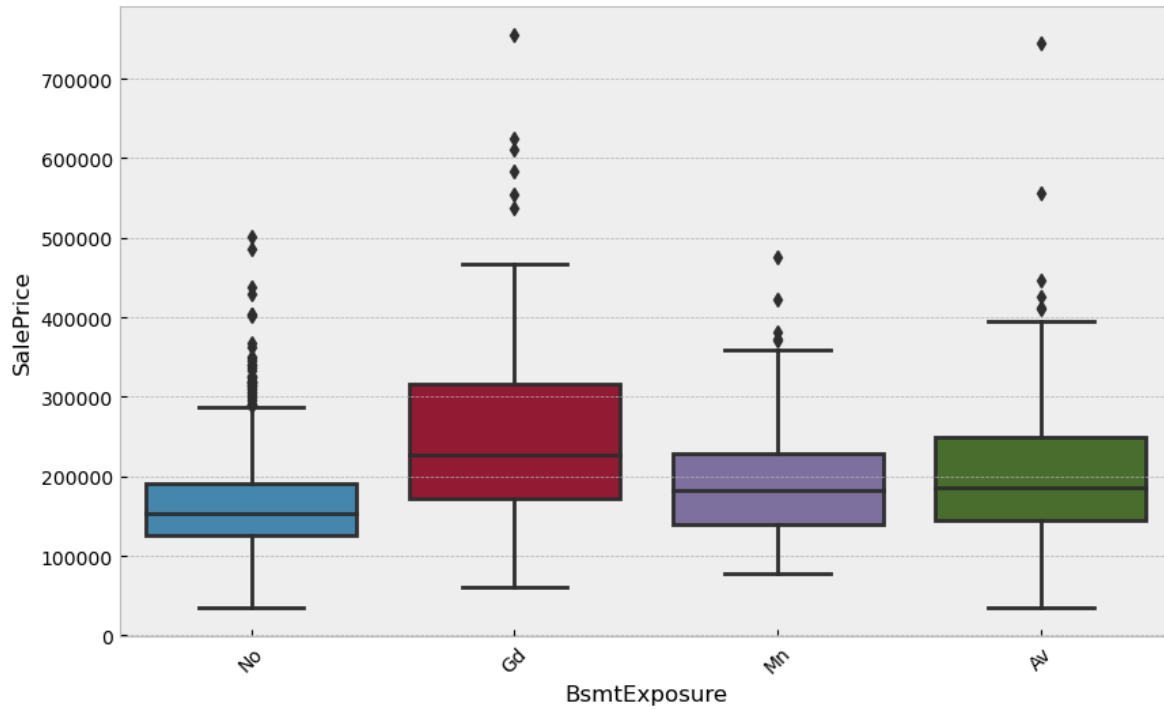
```
['MSZoning', 'Street', 'LotShape', 'LandContour', 'Utilities', 'LotConfig',
'LandSlope', 'Neighborhood', 'Condition1', 'Condition2', 'BldgType', 'HouseS
tyle', 'RoofStyle', 'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType',
'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual', 'BsmtCond', 'BsmtExposur
e', 'BsmtFinType1', 'BsmtFinType2', 'Heating', 'HeatingQC', 'CentralAir', 'E
lectrical', 'KitchenQual', 'Functional', 'FireplaceQu', 'GarageType', 'Garag
eFinish', 'GarageQual', 'GarageCond', 'PavedDrive', 'SaleType', 'SaleCondi
on']
```

In [21]:

```
plt.figure(figsize = (10, 6))  
ax = sns.boxplot(x='BsmtExposure', y='SalePrice', data=df_categ)  
plt.setp(ax.artists, alpha=.5, linewidth=2, edgecolor="k")  
plt.xticks(rotation=45)
```

Out[21]:

```
(array([0, 1, 2, 3]),  
 [Text(0, 0, 'No'), Text(1, 0, 'Gd'), Text(2, 0, 'Mn'), Text(3, 0, 'Av')])
```



In [22]:

```
plt.figure(figsize = (12, 6))
ax = sns.boxplot(x='SaleCondition', y='SalePrice', data=df_categ)
plt.setp(ax.artists, alpha=.5, linewidth=2, edgecolor="k")
plt.xticks(rotation=45)
```

Out[22]:

```
(array([0, 1, 2, 3, 4, 5]),
 [Text(0, 0, 'Normal'),
  Text(1, 0, 'Abnorml'),
  Text(2, 0, 'Partial'),
  Text(3, 0, 'AdjLand'),
  Text(4, 0, 'Alloca'),
  Text(5, 0, 'Family')])
```

