

1 House Prices: Advanced Regression Techniques

1.1 Description

Ask a home buyer to describe their dream house, and they probably won't begin with the height of the basement ceiling or the proximity to an east-west railroad. But this playground competition's dataset proves that much more influences price negotiations than the number of bedrooms or a white-picket fence.

With 79 explanatory variables describing (almost) every aspect of residential homes in Ames, Iowa, this competition challenges you to predict the final price of each home.

1.2 Acknowledgments

The Ames Housing dataset was compiled by Dean De Cock for use in data science education. It's an incredible alternative for data scientists looking for a modernized and expanded version of the often cited Boston Housing dataset.

1.3 Data set ¶

Here's a brief version of what you'll find in the data description file.

- **SalePrice** - the property's sale price in dollars. This is the target variable * that you're trying to predict.
- **MSSubClass** : The building class
- **MSZoning** : The general zoning classification
- **LotFrontage** : Linear feet of street connected to property
- **LotArea** : Lot size in square feet
- **Street** : Type of road access
- **Alley** : Type of alley access
- **LotShape** : General shape of property
- **LandContour** : Flatness of the property
- **Utilities** : Type of utilities available
- **LotConfig** : Lot configuration
- **LandSlope** : Slope of property
- **Neighborhood** : Physical locations within Ames city limits
- **Condition1** : Proximity to main road or railroad
- **Condition2** : Proximity to main road or railroad (if a second is present)
- **BldgType** : Type of dwelling
- **HouseStyle** : Style of dwelling
- **OverallQual** : Overall material and finish quality
- **OverallCond** : Overall condition rating
- **YearBuilt** : Original construction date
- **YearRemodAdd** : Remodel date
- **RoofStyle** : Type of roof
- **RoofMatl** : Roof material
- **Exterior1st** : Exterior covering on house
- **Exterior2nd** : Exterior covering on house (if more than one material)
- **MasVnrType** : Masonry veneer type
- **MasVnrArea** : Masonry veneer area in square feet
- **ExterQual** : Exterior material quality

- ExterCond : Present condition of the material on the exterior
- Foundation : Type of foundation
- BsmtQual : Height of the basement
- BsmtCond : General condition of the basement
- BsmtExposure : Walkout or garden level basement walls
- BsmtFinType1 : Quality of basement finished area
- BsmtFinSF1 : Type 1 finished square feet
- BsmtFinType2 : Quality of second finished area (if present)
- BsmtFinSF2 : Type 2 finished square feet
- BsmtUnfSF : Unfinished square feet of basement area
- TotalBsmtSF : Total square feet of basement area
- Heating : Type of heating
- HeatingQC : Heating quality and condition
- CentralAir : Central air conditioning
- Electrical : Electrical system
- 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
- BsmtFullBath : Basement full bathrooms
- BsmtHalfBath : Basement half bathrooms
- FullBath : Full bathrooms above grade
- HalfBath : Half baths above grade
- Bedroom : Number of bedrooms above basement level
- Kitchen : Number of kitchens
- KitchenQual : Kitchen quality
- TotRmsAbvGrd : Total rooms above grade (does not include bathrooms)
- Functional : Home functionality rating
- Fireplaces : Number of fireplaces
- FireplaceQu : Fireplace quality
- GarageType : Garage location
- GarageYrBlt : Year garage was built
- GarageFinish : Interior finish of the garage
- GarageCars : Size of garage in car capacity
- GarageArea : Size of garage in square feet
- GarageQual : Garage quality
- GarageCond : Garage condition
- PavedDrive : Paved driveway
- WoodDeckSF : Wood deck area in square feet
- OpenPorchSF : Open porch area in square feet
- EnclosedPorch : Enclosed porch area in square feet
- 3SsnPorch : Three season porch area in square feet
- ScreenPorch : Screen porch area in square feet
- PoolArea : Pool area in square feet
- PoolQC : Pool quality
- Fence : Fence quality
- MiscFeature : Miscellaneous feature not covered in other categories
- MiscVal : USD Value of miscellaneous feature
- MoSold : Month Sold
- YrSold : Year Sold
- SaleType : Type of sale
- SaleCondition : Condition of sale

1.4 Goal

Predict sales prices and practice feature engineering, RFs, and gradient boosting Type: supervised machine learning - regression

2 Exploratory Analysis¶

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

In [2]:

```
1 from sklearn.preprocessing import StandardScaler, MinMaxScaler
2 from sklearn.model_selection import train_test_split
```

In [3]:

```
1 # keep only relevant imports based on the regresssion or classification goals
2 from sklearn.metrics import mean_squared_error
3 from sklearn.model_selection import cross_val_score, GridSearchCV, RandomizedSe
```

In [4]:

```
1 # common classifiers
2 #from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
3 #from sklearn.linear_model import LogisticRegression, SGDClassifier
4 from sklearn.svm import SVC, LinearSVC
```

In [5]:

```
1 import xgboost as xgb
2 import lightgbm as lgbm
```

In [6]:

```
1 # common regresssors
2 from sklearn.linear_model import LinearRegression, Ridge, Lasso, ElasticNet, SG
3 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor,
4 from sklearn.svm import SVR
```

In [7]:

```
1 from sklearn.pipeline import Pipeline
```

In [8]:

```
1 # skip future warnings and display enough columns for wide data sets
2 import warnings
3 warnings.simplefilter(action='ignore') #, category=FutureWarning)
4 pd.set_option('display.max_columns', 100)
```

2.1 Data set first insight

Let's see with the data set looks like

In [9]:

```
1 df = pd.read_csv('../input/train.csv', index_col='Id' )
2 df.head()
```

Out[9]:

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

Number of samples (lines) and features (columns including the target)

In [10]:

```
1 df.shape
```

Out[10]:

(1460, 80)

Basic infos

In [11]:

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1460 entries, 1 to 1460
Data columns (total 80 columns):
MSSubClass      1460 non-null int64
MSZoning        1460 non-null object
LotFrontage     1201 non-null float64
LotArea         1460 non-null int64
Street          1460 non-null object
Alley           91 non-null object
LotShape        1460 non-null object
LandContour     1460 non-null object
Utilities       1460 non-null object
LotConfig       1460 non-null object
LandSlope       1460 non-null object
Neighborhood    1460 non-null object
Condition1      1460 non-null object
Condition2      1460 non-null object
BldgType        1460 non-null object
HouseStyle      1460 non-null object
OverallQual     1460 non-null int64
OverallCond     1460 non-null int64
YearBuilt       1460 non-null int64
YearRemodAdd    1460 non-null int64
RoofStyle       1460 non-null object
RoofMatl        1460 non-null object
Exterior1st     1460 non-null object
Exterior2nd     1460 non-null object
MasVnrType      1452 non-null object
MasVnrArea      1452 non-null float64
ExterQual       1460 non-null object
ExterCond       1460 non-null object
Foundation      1460 non-null object
BsmtQual        1423 non-null object
BsmtCond        1423 non-null object
BsmtExposure    1422 non-null object
BsmtFinType1    1423 non-null object
BsmtFinSF1      1460 non-null int64
BsmtFinType2    1422 non-null object
BsmtFinSF2      1460 non-null int64
BsmtUnfSF       1460 non-null int64
TotalBsmtSF     1460 non-null int64
Heating         1460 non-null object
HeatingQC       1460 non-null object
CentralAir      1460 non-null object
Electrical      1459 non-null object
1stFlrSF        1460 non-null int64
2ndFlrSF        1460 non-null int64
LowQualFinSF    1460 non-null int64
GrLivArea       1460 non-null int64
BsmtFullBath    1460 non-null int64
BsmtHalfBath    1460 non-null int64
FullBath        1460 non-null int64
HalfBath        1460 non-null int64
BedroomAbvGr    1460 non-null int64
KitchenAbvGr    1460 non-null int64
KitchenQual     1460 non-null object
TotRmsAbvGrd    1460 non-null int64
```

Functional	1460 non-null object
Fireplaces	1460 non-null int64
FireplaceQu	770 non-null object
GarageType	1379 non-null object
GarageYrBlt	1379 non-null float64
GarageFinish	1379 non-null object
GarageCars	1460 non-null int64
GarageArea	1460 non-null int64
GarageQual	1379 non-null object
GarageCond	1379 non-null object
PavedDrive	1460 non-null object
WoodDeckSF	1460 non-null int64
OpenPorchSF	1460 non-null int64
EnclosedPorch	1460 non-null int64
3SsnPorch	1460 non-null int64
ScreenPorch	1460 non-null int64
PoolArea	1460 non-null int64
PoolQC	7 non-null object
Fence	281 non-null object
MiscFeature	54 non-null object
MiscVal	1460 non-null int64
MoSold	1460 non-null int64
YrSold	1460 non-null int64
SaleType	1460 non-null object
SaleCondition	1460 non-null object
SalePrice	1460 non-null int64

dtypes: float64(3), int64(34), object(43)
memory usage: 923.9+ KB

Number of columns for each type of data

In [12]:

```
1 df.dtypes.value_counts()
```

Out[12]:

```
object      43
int64       34
float64      3
dtype: int64
```

Unique values for each type of data

In [13]:

```
1 df.select_dtypes('object').apply(pd.Series.nunique, axis = 0)
```

Out[13]:

MSZoning	5
Street	2
Alley	2
LotShape	4
LandContour	4
Utilities	2
LotConfig	5
LandSlope	3
Neighborhood	25
Condition1	9
Condition2	8
BldgType	5
HouseStyle	8
RoofStyle	6
RoofMatl	8
Exterior1st	15
Exterior2nd	16
MasVnrType	4
ExterQual	4
ExterCond	5
Foundation	6
BsmtQual	4
BsmtCond	4
BsmtExposure	4
BsmtFinType1	6
BsmtFinType2	6
Heating	6
HeatingQC	5
CentralAir	2
Electrical	5
KitchenQual	4
Functional	7
FireplaceQu	5
GarageType	6
GarageFinish	3
GarageQual	5
GarageCond	5
PavedDrive	3
PoolQC	3
Fence	4
MiscFeature	4
SaleType	9
SaleCondition	6
dtype:	int64

Ratio of missing values by column

In [14]:

```
1 def missing_values_table(df):
2     """Function to calculate missing values by column# Funct // credits Wil
3
4     # Total missing values
5     mis_val = df.isnull().sum()
6
7     # Percentage of missing values
8     mis_val_percent = 100 * df.isnull().sum() / len(df)
9
10    # Make a table with the results
11    mis_val_table = pd.concat([mis_val, mis_val_percent], axis=1)
12
13    # Rename the columns
14    mis_val_table_ren_columns = mis_val_table.rename(
15    columns = {0 : 'Missing Values', 1 : '% of Total Values'})
16
17    # Sort the table by percentage of missing descending
18    mis_val_table_ren_columns = mis_val_table_ren_columns[
19    mis_val_table_ren_columns.iloc[:,1] != 0].sort_values(
20    '% of Total Values', ascending=False).round(1)
21
22    # Print some summary information
23    print ("Le jeu de données a " + str(df.shape[1]) + " colonnes.\n"
24    "Il y a " + str(mis_val_table_ren_columns.shape[0]) +
25    " colonnes avec des valeurs manquantes.")
26
27    # Return the dataframe with missing information
28    return mis_val_table_ren_columns
```

In [15]:

```
1 missing_values = missing_values_table(df)
2 missing_values.head(10)
```

Le jeu de données a 80 colonnes.
Il y a 19 colonnes avec des valeurs manquantes.

Out[15]:

	Missing Values	% of Total Values
PoolQC	1453	99.5
MiscFeature	1406	96.3
Alley	1369	93.8
Fence	1179	80.8
FireplaceQu	690	47.3
LotFrontage	259	17.7
GarageType	81	5.5
GarageYrBlt	81	5.5
GarageFinish	81	5.5
GarageQual	81	5.5

In [16]:

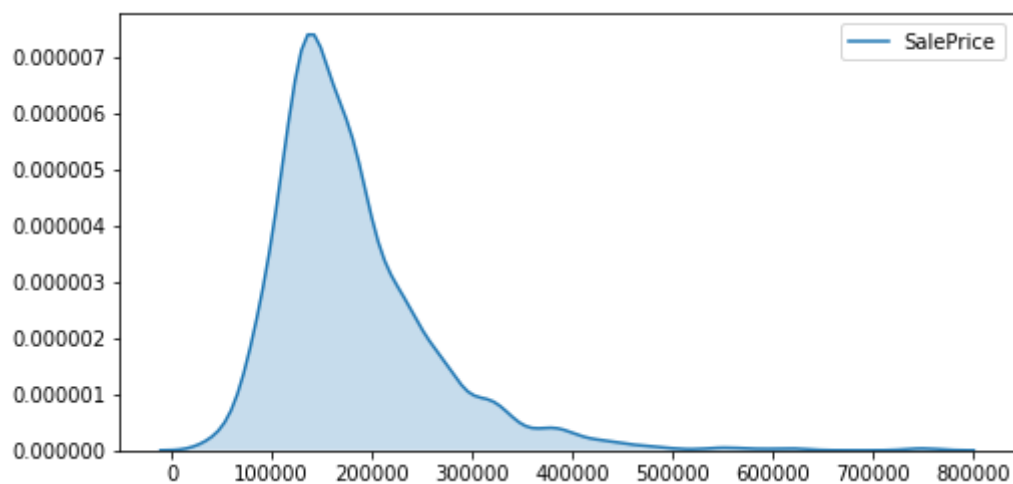
```
1 cat_feat = list(df.select_dtypes('object').columns)
2 num_feat = list(df.select_dtypes(exclude='object').columns)
```

2.2 Data Visualization

informations on the target

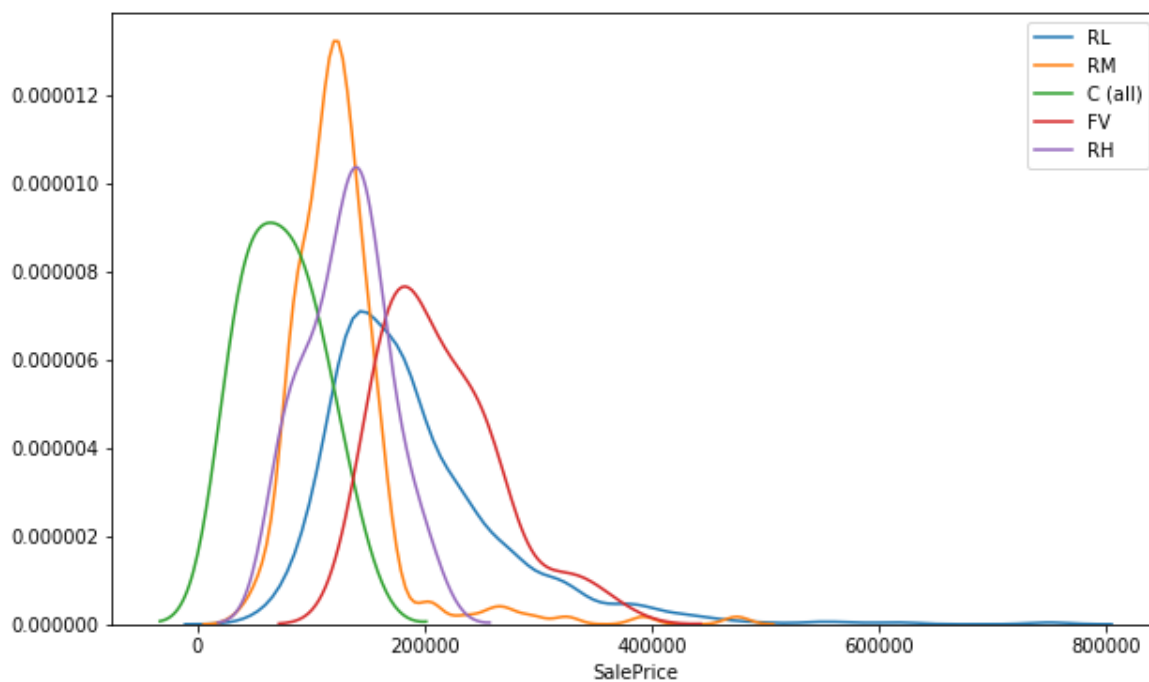
In [17]:

```
1 plt.figure(figsize=(8, 4))
2 sns.kdeplot(df.SalePrice, shade=True)
3 plt.show()
```



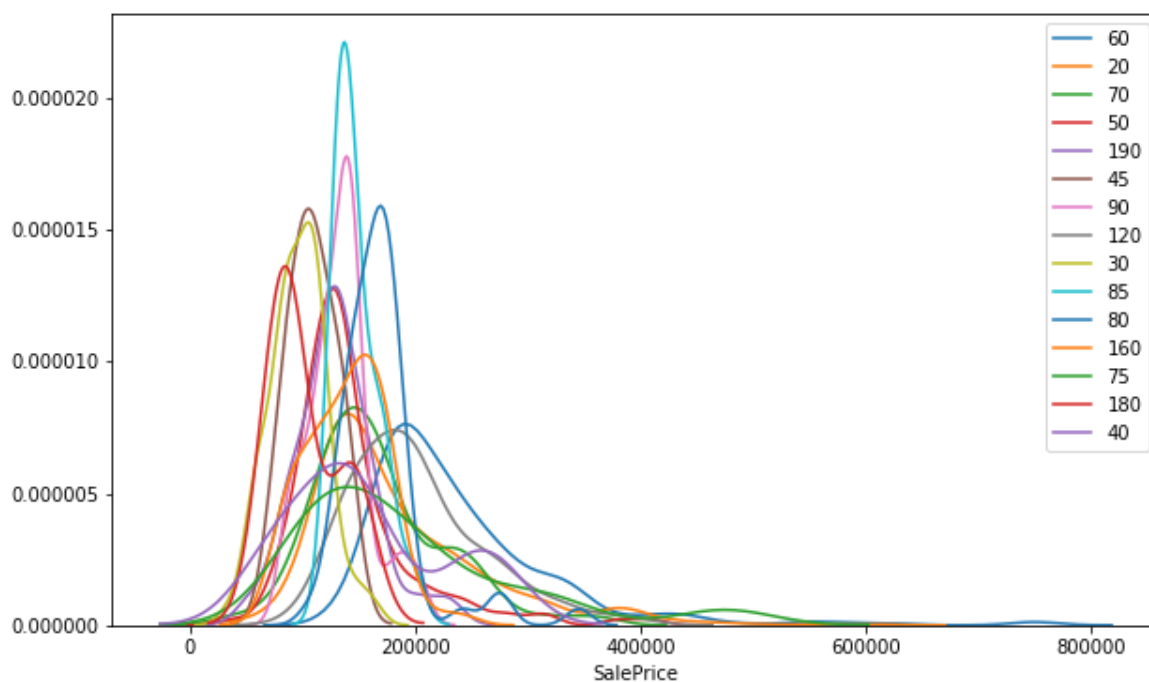
In [18]:

```
1 plt.figure(figsize=(10, 6))
2 for zone in list(df.MSZoning.unique()):
3     sns.distplot(df[df.MSZoning==zone].SalePrice, label=zone, hist=False)
4 plt.show()
```



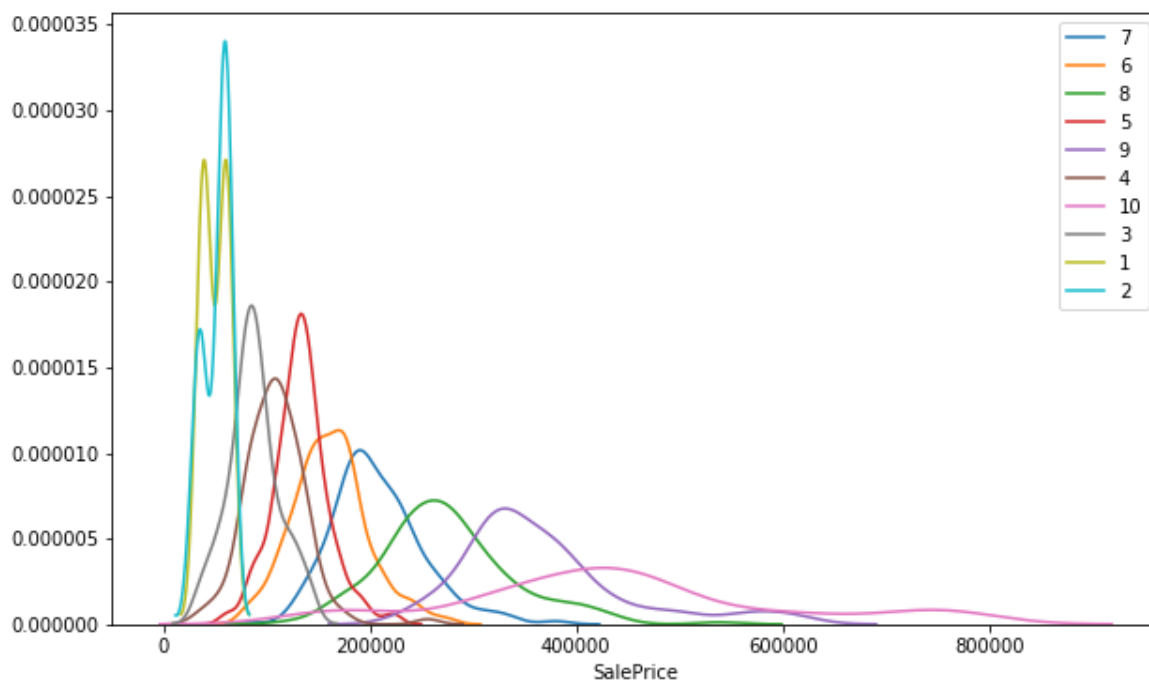
In [19]:

```
1 plt.figure(figsize=(10, 6))
2 for ms_sub_class in list(df.MSSubClass.unique()):
3     sns.distplot(df[df.MSSubClass==ms_sub_class].SalePrice, label=ms_sub_class,
4 plt.show()
5
```



In [20]:

```
1 plt.figure(figsize=(10, 6))
2 for qual in list(df.OverallQual.unique()):
3     sns.distplot(df[df.OverallQual==qual].SalePrice, label=qual, hist=False)
4 plt.show()
```



In [21]:

```
1 df.SalePrice.describe()
```

Out[21]:

```
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
```

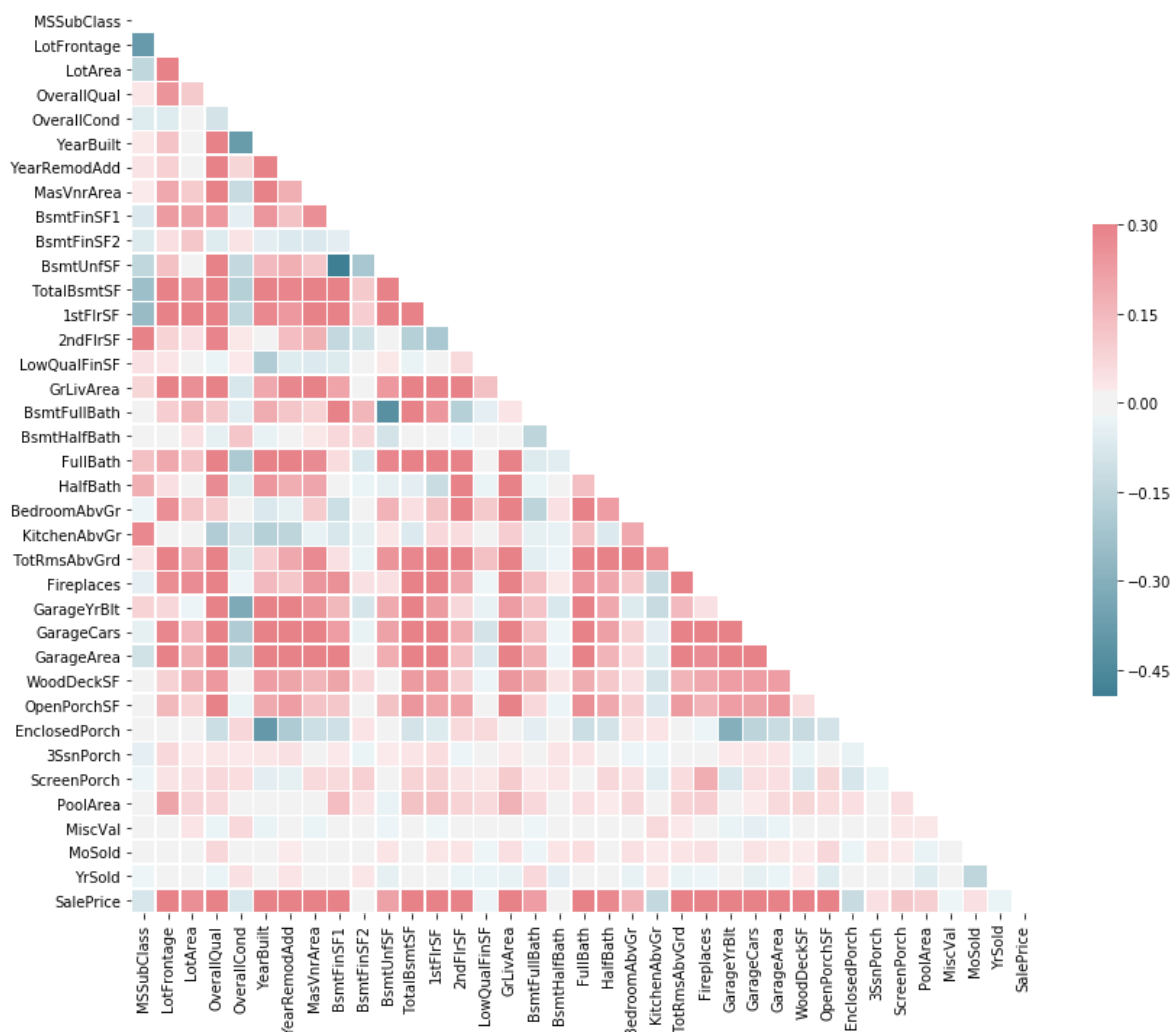
2.3 Correlations

In [22]:

```
1 corr = df.corr()
2 corr
3
4 # Generate a mask for the upper triangle
5 mask = np.zeros_like(corr, dtype=np.bool)
6 mask[np.triu_indices_from(mask)] = True
7
8 # Set up the matplotlib figure
9 f, ax = plt.subplots(figsize=(14, 12))
10
11 # Generate a custom diverging colormap
12 cmap = sns.diverging_palette(220, 10, as_cmap=True)
13
14 # Draw the heatmap with the mask and correct aspect ratio
15 sns.heatmap(corr, mask=mask, cmap=cmap, vmax=.3, center=0, square=True, linewidth
```

Out[22]:

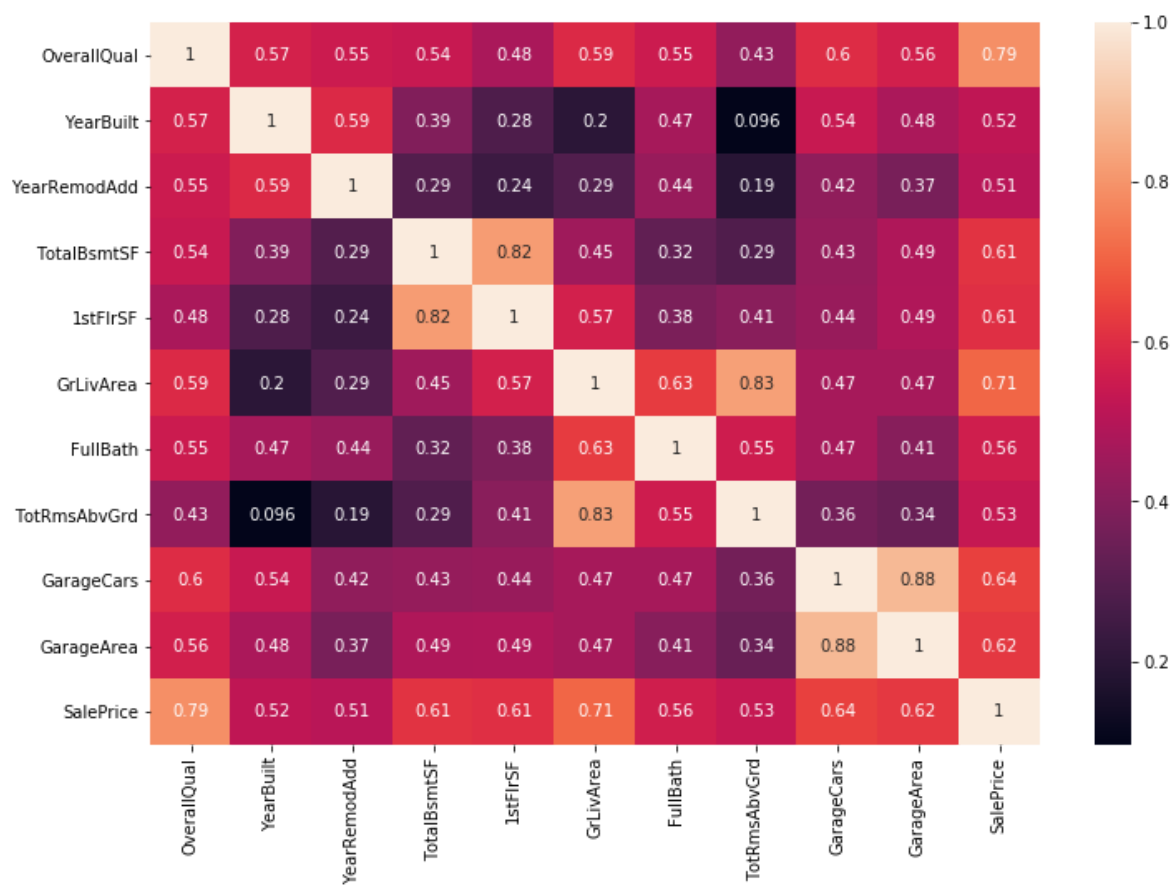
<matplotlib.axes._subplots.AxesSubplot at 0x7fbbf318cfd0>



Top 50% Corralation train attributes with sale-price

In [23]:

```
1 top_feature = corr.index[abs(corr['SalePrice'])>0.5]
2 plt.subplots(figsize=(12, 8))
3 top_corr = df[top_feature].corr()
4 sns.heatmap(top_corr, annot=True)
5 plt.show()
```



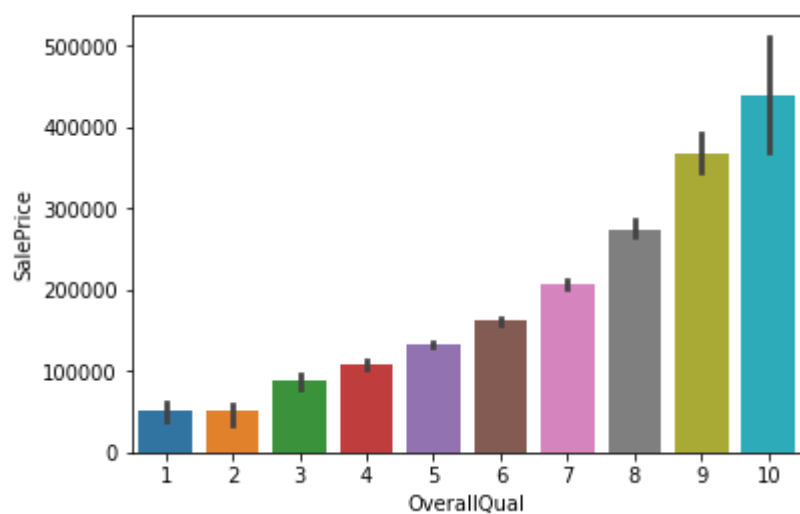
OverallQual is highly correlated with target feature of saleprice by near 80%

In [24]:

```
1 sns.barplot(df.OverallQual, df.SalePrice)
```

Out[24]:

<matplotlib.axes._subplots.AxesSubplot at 0x7fbbf16a0550>

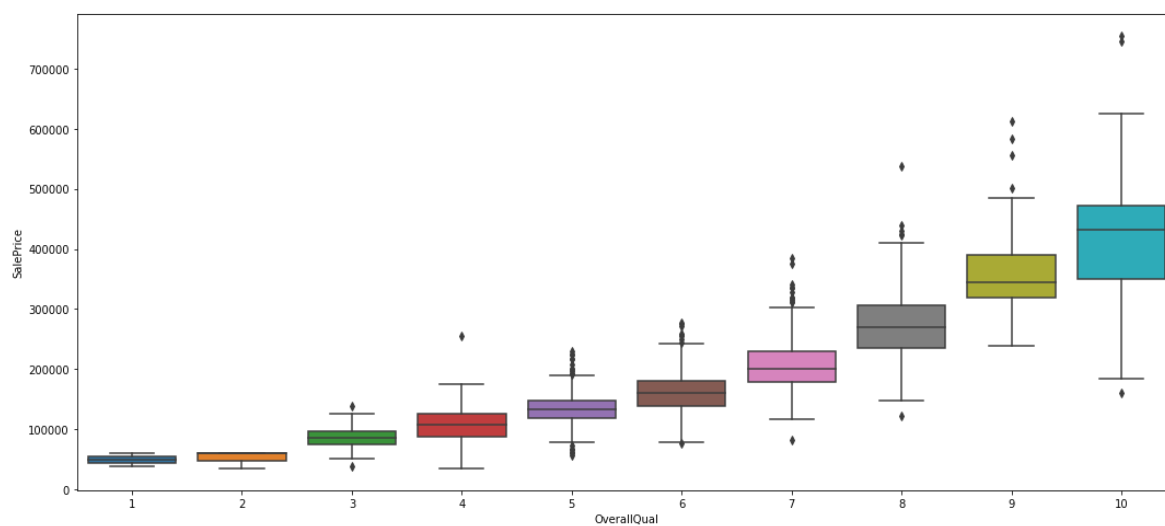


In [25]:

```
1 plt.figure(figsize=(18, 8))
2 sns.boxplot(x=df.OverallQual, y=df.SalePrice)
```

Out[25]:

<matplotlib.axes._subplots.AxesSubplot at 0x7fbbf1335a20>

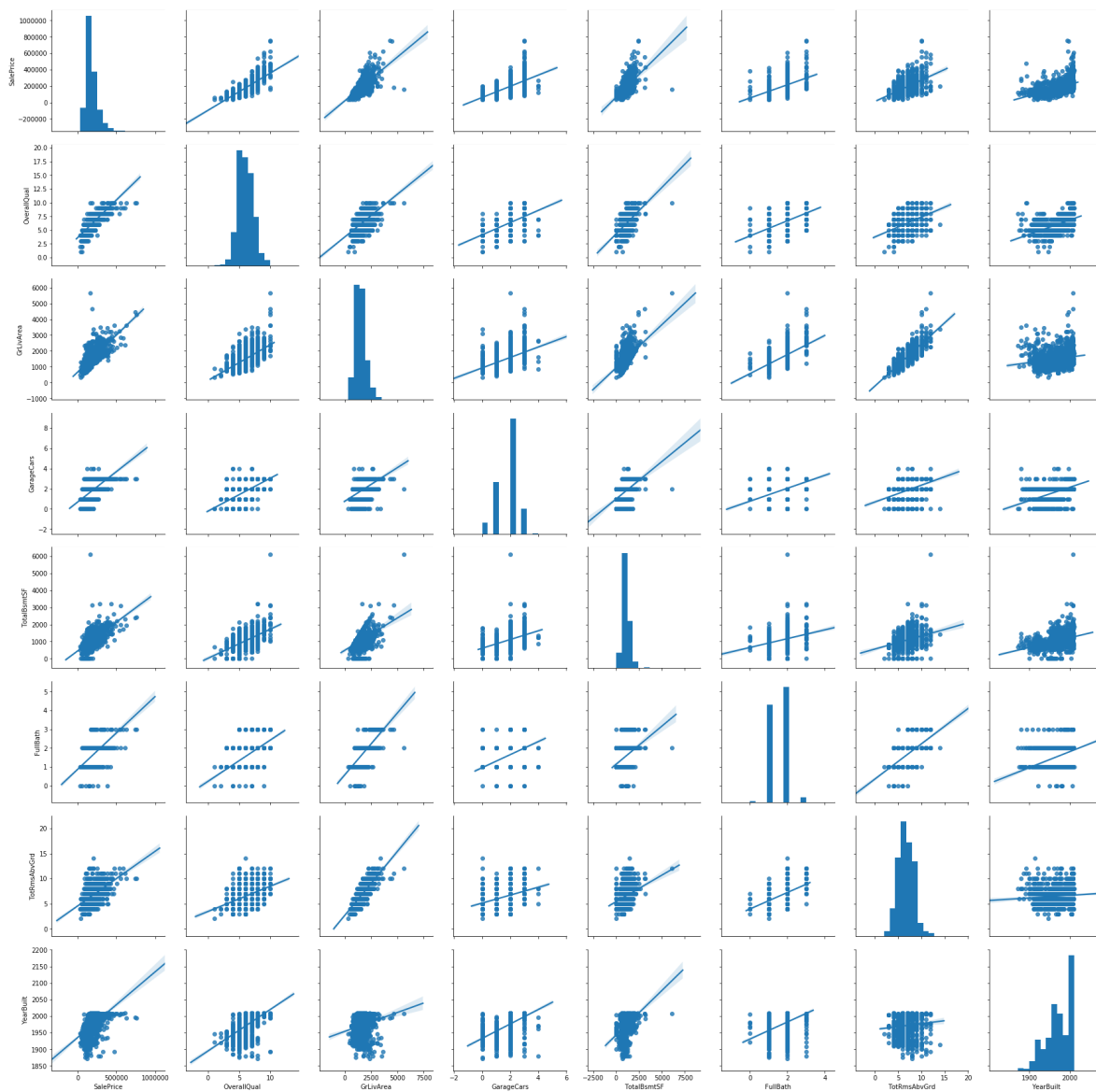


In [26]:

```
1 col = ['SalePrice', 'OverallQual', 'GrLivArea', 'GarageCars', 'TotalBsmtSF', 'F  
2 sns.pairplot(df[col], height=3, kind='reg')
```

Out[26]:

<seaborn.axisgrid.PairGrid at 0x7fbbf1668898>



In [27]:

```
1 print("Most postively correlated features with the target")
2 corr = df.corr()
3 corr.sort_values(['SalePrice'], ascending=False, inplace=True)
4 corr.SalePrice
```

Most postively correlated features with the target

Out[27]:

SalePrice	1.000000
OverallQual	0.790982
GrLivArea	0.708624
GarageCars	0.640409
GarageArea	0.623431
TotalBsmstSF	0.613581
1stFlrSF	0.605852
FullBath	0.560664
TotRmsAbvGrd	0.533723
YearBuilt	0.522897
YearRemodAdd	0.507101
GarageYrBltd	0.486362
MasVnrArea	0.477493
Fireplaces	0.466929
BsmstFinSF1	0.386420
LotFrontage	0.351799
WoodDeckSF	0.324413
2ndFlrSF	0.319334
OpenPorchSF	0.315856
HalfBath	0.284108
LotArea	0.263843
BsmstFullBath	0.227122
BsmstUnfSF	0.214479
BedroomAbvGr	0.168213
ScreenPorch	0.111447
PoolArea	0.092404
MoSold	0.046432
3SsnPorch	0.044584
BsmstFinSF2	-0.011378
BsmstHalfBath	-0.016844
MiscVal	-0.021190
LowQualFinSF	-0.025606
YrSold	-0.028923
OverallCond	-0.077856
MSSubClass	-0.084284
EnclosedPorch	-0.128578
KitchenAbvGr	-0.135907

Name: SalePrice, dtype: float64

3 Data preparation & feature engineering

3.1 Dealing with abnormal values

Not relevant here, we can assume that all values are been well integrated.

3.2 Data cleaning & Label encoding of categorical features

No duplicated rows

In [28]:

```
1 df.duplicated().sum()
```

Out[28]:

0

Let's remove columns with a high ratio of missing values

We don't have much samples, so instead of removing rows with nan, missing values are then replaced by the median

In [29]:

```
1 from sklearn.preprocessing import LabelEncoder
```

In [30]:

```
1 def prepare_data(dataframe):
2
3     dataframe = dataframe.drop(columns=['PoolQC', 'MiscFeature', 'Alley', 'Fenc
4
5     cat_feat = list(dataframe.select_dtypes('object').columns)
6     num_feat = list(dataframe.select_dtypes(exclude='object').columns)
7
8     dataframe[num_feat] = dataframe[num_feat].fillna(dataframe[num_feat].median
9     dataframe[cat_feat] = dataframe[cat_feat].fillna("Not communicated")
10
11     for c in cat_feat:
12         lbl = LabelEncoder()
13         lbl.fit(list(dataframe[c].values))
14         dataframe[c] = lbl.transform(list(dataframe[c].values))
15
16     return dataframe
```

At first sight, there isn't any value in the wrong type / format

Those features can't be used as they are (in string format), this is why we need to convert them in a numerical way...

In [31]:

```
1 df = prepare_data(df)
```

3.3 Creation of new features

- In this case, it's complicated to add features from an other dataset because no information is provided with the CSV file we're using.
- All columns except the id (used as index) seems to be relevant, so all of them are kept at first.
- We can also combine features to create new ones - but in this case it doesn't seem to be really usefull.

3.4 Standardization / normalization

Not needed here

In [32]:

```
1 #df[num_feat] = MinMaxScaler().fit_transform(df[num_feat])
```

3.5 Feature selection & and data preparation for models

In [33]:

```
1 y = df['SalePrice']
2 X = df.drop(columns=['SalePrice'])
3 X.shape, y.shape
```

Out[33]:

```
((1460, 74), (1460,))
```

Let's split the data into a train and a test set

In [34]:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
2 X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

Out[34]:

```
((1168, 74), (292, 74), (1168,), (292,))
```

3.6 Feature importance

Top 10 most important features:

In [35]:

```
1 rnd_reg = RandomForestRegressor(n_estimators=500, n_jobs=-1)
2 rnd_reg.fit(X, y)
3
4 feature_importances = pd.DataFrame(rnd_reg.feature_importances_, index = X.columns,
5                                   columns=['importance']).sort_values('importance')
```

In [36]:

```
1 feature_importances[:10]
```

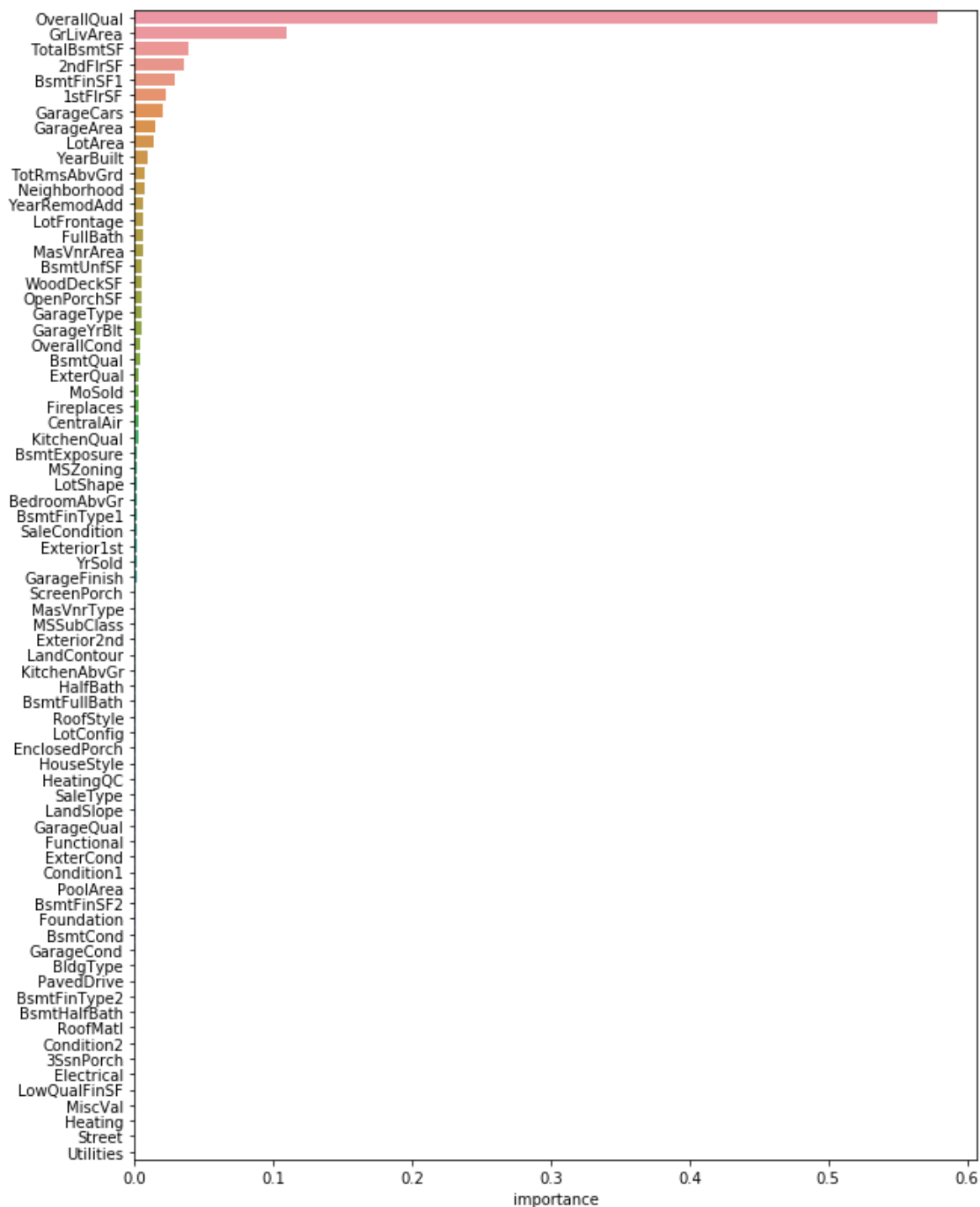
Out[36]:

	importance
OverallQual	0.578202
GrLivArea	0.109874
TotalBsmtSF	0.039231
2ndFlrSF	0.035708
BsmtFinSF1	0.029386
1stFlrSF	0.022536
GarageCars	0.020744
GarageArea	0.015562
LotArea	0.013600
YearBuilt	0.009355

Graph with features sorted by importance

In [37]:

```
1 plt.figure(figsize=(10, 14))
2 sns.barplot(x="importance", y=feature_importances.index, data=feature_importanc
3 plt.show()
```



4 Training models and results

4.1 Baselines - first selection of models

In [38]:

```
1 # f1_score binary by default
2 def get_rmse(reg, model_name):
3     """Print the score for the model passed in argument and retrun scores for t
4
5     y_train_pred, y_pred = reg.predict(X_train), reg.predict(X_test)
6     rmse_train, rmse_test = np.sqrt(mean_squared_error(y_train, y_train_pred)),
7     print(model_name, f'\t - RMSE on Training = {rmse_train:.0f} / RMSE on Tes
8
9     return rmse_train, rmse_test
```

In [39]:

```
1 model_list = [
2     LinearRegression(), Lasso(), SVR(),
3     RandomForestRegressor(), GradientBoostingRegressor(), Ridge(), ElasticNet()
4     BayesianRidge(), ExtraTreesRegressor()
5 ]
```

In [40]:

```
1 model_names = [str(m)[:str(m).index('(')] for m in model_list]
2 rmse_train, rmse_test = [], []
```

In [41]:

```
1 model_names
```

Out[41]:

```
['LinearRegression',
 'Lasso',
 'SVR',
 'RandomForestRegressor',
 'GradientBoostingRegressor',
 'Ridge',
 'ElasticNet',
 'LinearSVC',
 'BayesianRidge',
 'ExtraTreesRegressor']
```

In [42]:

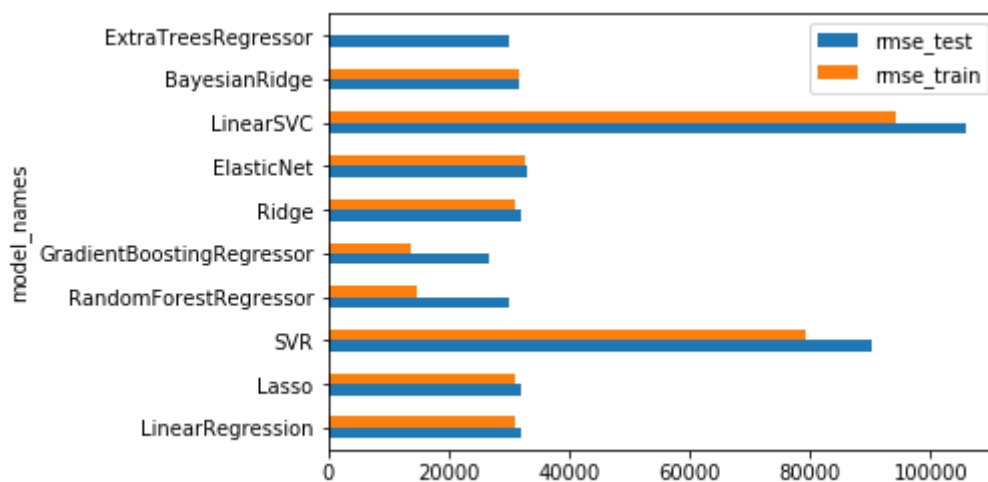
```
1 for model, name in zip(model_list, model_names):
2     model.fit(X_train, y_train)
3     sc_train, sc_test = get_rmse(model, name)
4     rmse_train.append(sc_train)
5     rmse_test.append(sc_test)
```

```
LinearRegression      - RMSE on Training  = 31163 / RMSE on Test =
32162
Lasso      - RMSE on Training  = 31163 / RMSE on Test = 32158
SVR        - RMSE on Training  = 79338 / RMSE on Test = 90251
RandomForestRegressor - RMSE on Training  = 14748 / RMSE on Test =
30100
GradientBoostingRegressor - RMSE on Training  = 13689 / RMSE on
Test = 26783
Ridge      - RMSE on Training  = 31176 / RMSE on Test = 32091
ElasticNet  - RMSE on Training  = 32547 / RMSE on Test = 33122
LinearSVC   - RMSE on Training  = 94350 / RMSE on Test = 105986
BayesianRidge - RMSE on Training  = 31599 / RMSE on Test = 31864
ExtraTreesRegressor - RMSE on Training  = 0 / RMSE on Test = 3002
3
```

Results comparison chart

In [43]:

```
1 df_score = pd.DataFrame({'model_names' : model_names,
2                           'rmse_train' : rmse_train,
3                           'rmse_test' : rmse_test})
4 ax = df_score.plot.barh(y=['rmse_test', 'rmse_train'], x='model_names')
```



The LinearSVC model isn't performing well because data haven't been scaled before, let's do it with a pipeline:

In [44]:

```
1 svm_reg = Pipeline([
2     ("scaler", StandardScaler()),
3     ("svm_regressor", LinearSVC())
4 ])
5 svm_reg.fit(X_train, y_train)
6 _, _ = get_rmse(svm_reg, "svr_rbf")
```

svr_rbf - RMSE on Training = 2158 / RMSE on Test = 70136

That's much better, although it seems the linear kernel is the best option here:

In [45]:

```
1 svr_rbf = SVR(kernel = 'rbf')
2 svr_rbf.fit(X_train, y_train)
3 _, _ = get_rmse(svr_rbf, "svr_rbf")
```

svr_rbf - RMSE on Training = 79338 / RMSE on Test = 90251

In [46]:

```
1 svm_reg = Pipeline([
2     ("scaler", StandardScaler()),
3     ("svm_regressor", SVR())
4 ])
5 svm_reg.fit(X_train, y_train)
6 _, _ = get_rmse(svm_reg, "svr_rbf")
7
8 svm_reg = Pipeline([
9     ("scaler", StandardScaler()),
10    ("svm_regressor", SVR(kernel="poly"))
11 ])
12 svm_reg.fit(X_train, y_train)
13 _, _ = get_rmse(svm_reg, "svr_poly")
14
15 sgd_reg = Pipeline([
16     ("scaler", StandardScaler()),
17     ("sgd_regressor", SGDRegressor())
18 ])
19 sgd_reg.fit(X_train, y_train)
20 _, _ = get_rmse(sgd_reg, "sgd_reg")
```

svr_rbf - RMSE on Training = 79310 / RMSE on Test = 90222
svr_poly - RMSE on Training = 79315 / RMSE on Test = 90224
sgd_reg - RMSE on Training = 31546 / RMSE on Test = 34055

The same remark comes true also for the SGD Regressor model

Let's try XGBoost !

In [47]:

```
1 xgb_reg = xgb.XGBRegressor()
2 xgb_reg.fit(X_train, y_train)
3 _, _ = get_rmse(xgb_reg, "xgb_reg")
```

xgb_reg - RMSE on Training = 14601 / RMSE on Test = 27209

Looks promising, here we can conclude that RandomForestRegressor, GradientBoostingRegressor and XGBoost seems to be the models we'll keep for hyperparameters tuning !

4.2 Model optimisation

4.2.1 RandomForrestReg

In [48]:

```
1 from sklearn.model_selection import GridSearchCV
2
3
4 rf = RandomForestRegressor()
5 param_grid = {
6     'n_estimators': [80, 100, 120],
7     'max_features': [14, 15, 16, 17],
8     'max_depth' : [14, 16, 18]
9 }
10
11
12 rfc_cv = GridSearchCV(estimator=rf, param_grid=param_grid, cv=5, n_jobs=-1)
13 rfc_cv.fit(X_train, y_train)
14 print(rfc_cv.best_params_)
15 _, _ = get_rmse(rfc_cv, "rfc_reg")
```

{'max_depth': 18, 'max_features': 17, 'n_estimators': 100}
rfc_reg - RMSE on Training = 11404 / RMSE on Test = 29079

4.2.2 GradientBoostingReg

In [49]:

```
1 gb = GradientBoostingRegressor()
2 param_grid = {
3     'n_estimators': [100, 400],
4     'max_features': [14, 15, 16, 17],
5     'max_depth' : [1, 2, 8, 14, 18]
6 }
7
8
9 gb_cv = GridSearchCV(estimator=gb, param_grid=param_grid, cv=5, n_jobs=-1)
10 gb_cv.fit(X_train, y_train)
11 print(gb_cv.best_params_)
12 _, _ = get_rmse(gb_cv, "gb_cv")
```

{'max_depth': 8, 'max_features': 15, 'n_estimators': 100}
gb_cv - RMSE on Training = 1180 / RMSE on Test = 25624

4.2.3 XGBoostReg

In [50]:

```
1 xg = xgb.XGBRegressor()
2 param_grid = {
3     'n_estimators': [100, 400],
4     'max_features': [10, 14, 16],
5     'max_depth' : [1, 2, 8, 18]
6 }
7
8
9 xg_cv = GridSearchCV(estimator=xg, param_grid=param_grid, cv=5, n_jobs=-1)
10 xg_cv.fit(X_train, y_train)
11 print(xg_cv.best_params_)
12 _, _ = get_rmse(xg_cv, "xg_cv")
```

```
{'max_depth': 8, 'max_features': 10, 'n_estimators': 100}
xg_cv      - RMSE on Training  = 2478 / RMSE on Test = 28332
```

4.3 Combination of the best models & submission

In [51]:

```
1 df_test = pd.read_csv('../input/test.csv', index_col='Id' )
2 df_test.head()
```

Out[51]:

	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Ut
Id									
1461	20	RH	80.0	11622	Pave	NaN	Reg	Lvl	/
1462	20	RL	81.0	14267	Pave	NaN	IR1	Lvl	/
1463	60	RL	74.0	13830	Pave	NaN	IR1	Lvl	/
1464	60	RL	78.0	9978	Pave	NaN	IR1	Lvl	/
1465	120	RL	43.0	5005	Pave	NaN	IR1	HLS	/

In [52]:

```
1 df_test = prepare_data(df_test)
2 df_test.shape
```

Out[52]:

(1459, 74)

In [53]:

```
1 rfc_sub, gb_sub, xg_sub = rfc_cv.predict(df_test), gb_cv.predict(df_test), xg_c
```

In [54]:

```
1 sub = pd.DataFrame()  
2 sub['Id'] = df_test.index  
3 sub['SalePrice'] = np.mean([rfc_sub, gb_sub, xg_sub], axis=0) / 3  
4 sub.to_csv('submission.csv', index=False)
```

If you found this notebook helpful or you just liked it , some upvotes would be very much appreciated - That will keep me motivated to update it on a regular basis :-)