

# How I made top 0.3% on a Kaggle competition

Getting started with competitive data science can be quite intimidating. So I wrote this quick overview of how I made top 0.3% on the Advanced Regression Techniques competition. If there is interest, I'm happy to do deep dives into the intuition behind the feature engineering and models used in this kernel.

I encourage you to fork this kernel, play with the code and enter the competition. Good luck!

If you like this kernel, please give it an upvote. Thank you!

## The Goal

- Each row in the dataset describes the characteristics of a house.
- Our goal is to predict the SalePrice, given these features.
- Our models are evaluated on the Root-Mean-Squared-Error (RMSE) between the log of the SalePrice predicted by our model, and the log of the actual SalePrice. Converting RMSE errors to a log scale ensures that errors in predicting expensive houses and cheap houses will affect our score equally.

## Key features of the model training process in this kernel:

- **Cross Validation:** Using 12-fold cross-validation
- **Models:** On each run of cross-validation I fit 7 models (ridge, svr, gradient boosting, random forest, xgboost, lightgbm regressors)
- **Stacking:** In addition, I trained a meta StackingCVRegressor optimized using xgboost
- **Blending:** All models trained will overfit the training data to varying degrees. Therefore, to make final predictions, I blended their predictions together to get more robust predictions.

## Model Performance

We can observe from the graph below that the blended model far outperforms the other models, with an RMSLE of 0.075. This is the model I used for making the final predictions.

```
In [1]: from IPython.display import Image
        Image("../input/kernel-files/model_training_advanced_regression.png")
```

Out[1]:



Now that we have some context, let's get started!

In [2]:

```
# Essentials
import numpy as np
import pandas as pd
import datetime
import random

# Plots
import seaborn as sns
import matplotlib.pyplot as plt

# Models
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor, AdaBoostRegressor, BaggingRegressor
from sklearn.kernel_ridge import KernelRidge
from sklearn.linear_model import Ridge, RidgeCV
from sklearn.linear_model import ElasticNet, ElasticNetCV
from sklearn.svm import SVR
from mlxtend.regressor import StackingCVRegressor
import lightgbm as lgb
from lightgbm import LGBMRegressor
from xgboost import XGBRegressor

# Stats
from scipy.stats import skew, norm
from scipy.special import boxcox1p
from scipy.stats import boxcox_normmax

# Misc
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import KFold, cross_val_score
from sklearn.metrics import mean_squared_error
from sklearn.preprocessing import OneHotEncoder
from sklearn.preprocessing import LabelEncoder
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import scale
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import RobustScaler
from sklearn.decomposition import PCA

pd.set_option('display.max_columns', None)

# Ignore useless warnings
import warnings
warnings.filterwarnings(action="ignore")
pd.options.display.max_seq_items = 8000
pd.options.display.max_rows = 8000

import os
print(os.listdir("../input/kernel-files"))

['model_training_advanced_regression.png']
```

In [3]:

```
# Read in the dataset as a dataframe
train = pd.read_csv('../input/house-prices-advanced-regression-techniques/train.csv')
test = pd.read_csv('../input/house-prices-advanced-regression-techniques/test.csv')
```

```
train.shape, test.shape
```

```
Out[3]:  
((1460, 81), (1459, 80))
```

## EDA

### The Goal

- Each row in the dataset describes the characteristics of a house.
- Our goal is to predict the SalePrice, given these features.

```
In [4]:  
# Preview the data we're working with  
train.head()
```

```
Out[4]:
```

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

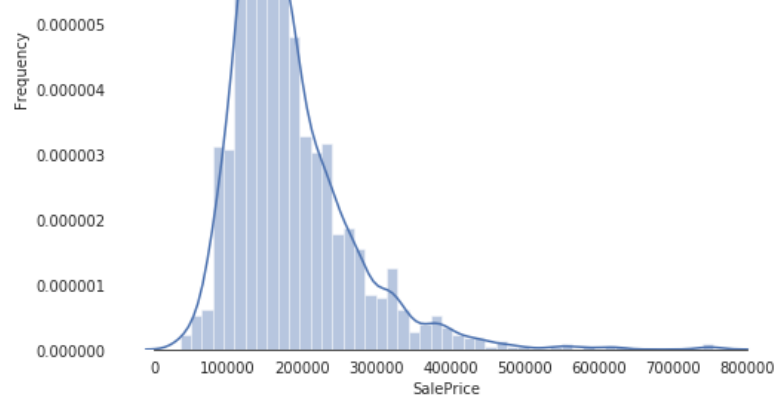
### SalePrice: the variable we're trying to predict

```
In [5]:  
sns.set_style("white")  
sns.set_color_codes(palette='deep')  
f, ax = plt.subplots(figsize=(8, 7))  
#Check the new distribution  
sns.distplot(train['SalePrice'], color="b");  
ax.xaxis.grid(False)  
ax.set(ylabel="Frequency")  
ax.set(xlabel="SalePrice")  
ax.set(title="SalePrice distribution")  
sns.despine(trim=True, left=True)  
plt.show()
```

SalePrice distribution

0.000008  
0.000007  
0.000006





In [6]:

```
# Skew and kurt
print("Skewness: %f" % train['SalePrice'].skew())
print("Kurtosis: %f" % train['SalePrice'].kurt())
```

Skewness: 1.882876

Kurtosis: 6.536282

## Features: a deep dive

Let's visualize some of the features in the dataset

In [7]:

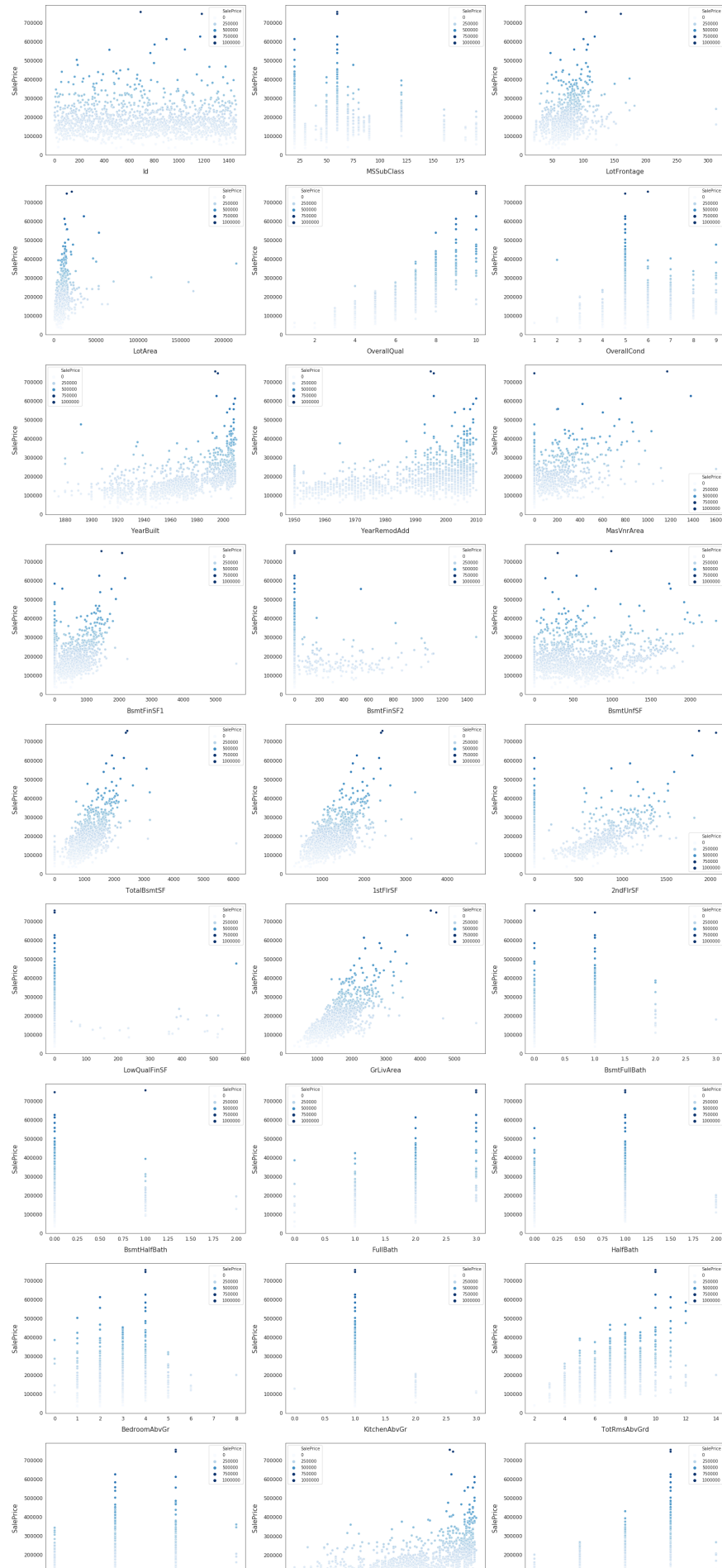
```
# Finding numeric features
numeric_dtypes = ['int16', 'int32', 'int64', 'float16', 'float32',
                  'float64']
numeric = []
for i in train.columns:
    if train[i].dtype in numeric_dtypes:
        if i in ['TotalSF', 'Total_Bathrooms', 'Total_porch_sf', 'haspool',
                'hasgarage', 'hasbsmt', 'hasfireplace']:
            pass
        else:
            numeric.append(i)
# visualising some more outliers in the data values
fig, axs = plt.subplots(ncols=2, nrows=0, figsize=(12, 120))
plt.subplots_adjust(right=2)
plt.subplots_adjust(top=2)
sns.color_palette("husl", 8)
for i, feature in enumerate(list(train[numeric]), 1):
    if(feature=='MiscVal'):
        break
    plt.subplot(len(list(numeric)), 3, i)
    sns.scatterplot(x=feature, y='SalePrice', hue='SalePrice', palette='Blues', data=train)

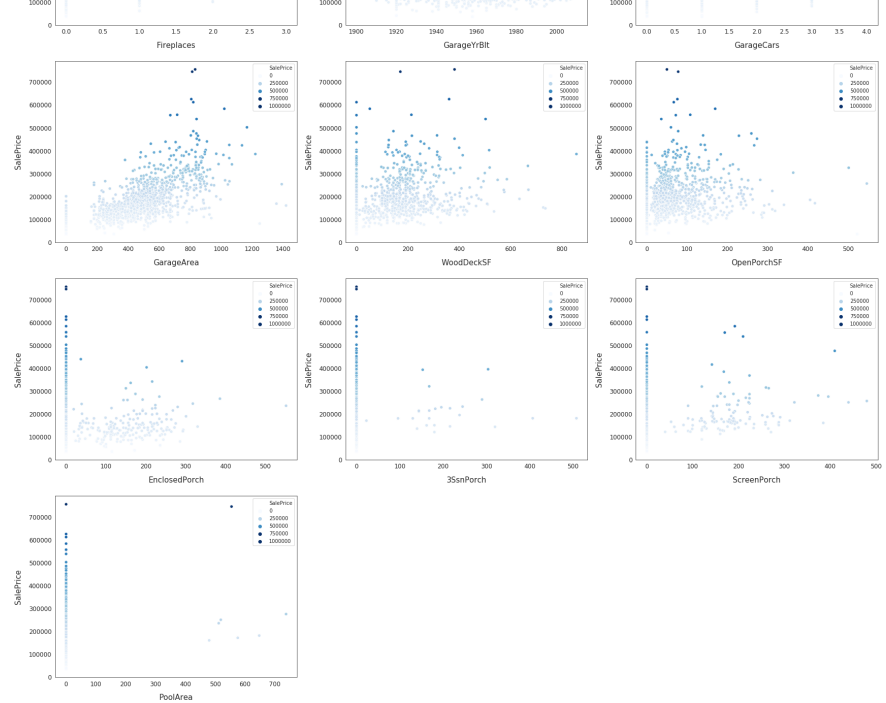
    plt.xlabel('{}'.format(feature), size=15, labelpad=12.5)
    plt.ylabel('SalePrice', size=15, labelpad=12.5)

    for j in range(2):
        plt.tick_params(axis='x', labelsize=12)
        plt.tick_params(axis='y', labelsize=12)
```

```
plt.legend(loc='best', prop={'size': 10})
```

```
plt.show()
```





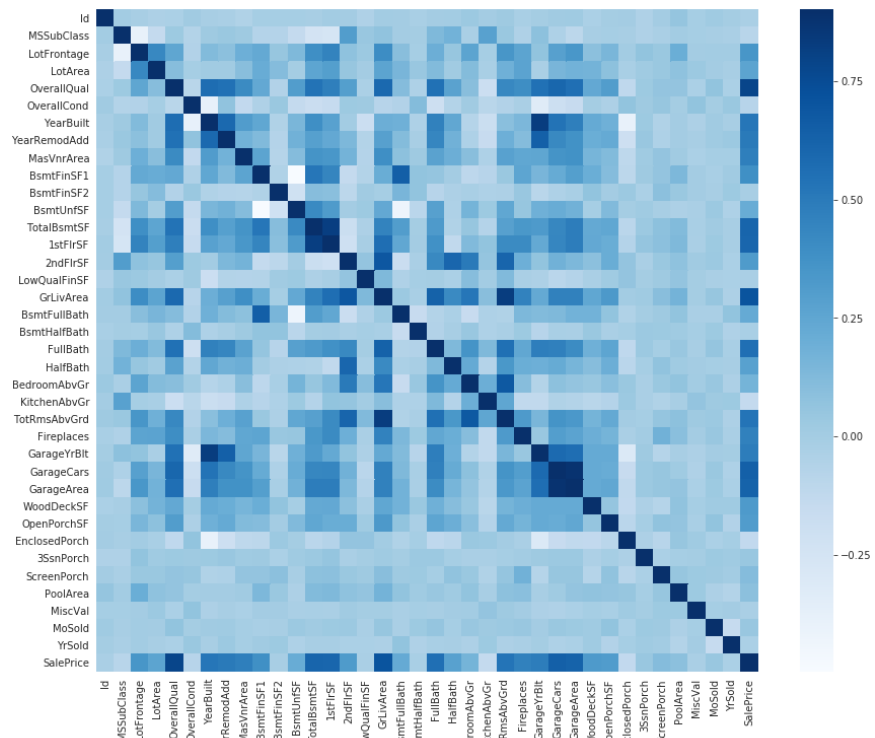
and plot how the features are correlated to each other, and to SalePrice

In [8]:

```
corr = train.corr()
plt.subplots(figsize=(15,12))
sns.heatmap(corr, vmax=0.9, cmap="Blues", square=True)
```

Out[8]:

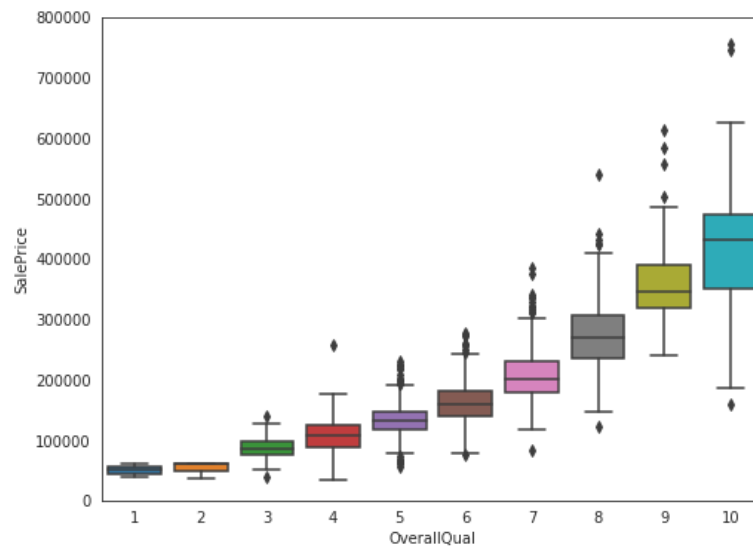
<matplotlib.axes.\_subplots.AxesSubplot at 0x7ff0e416e4e0>



Let's plot how SalePrice relates to some of the features in the dataset

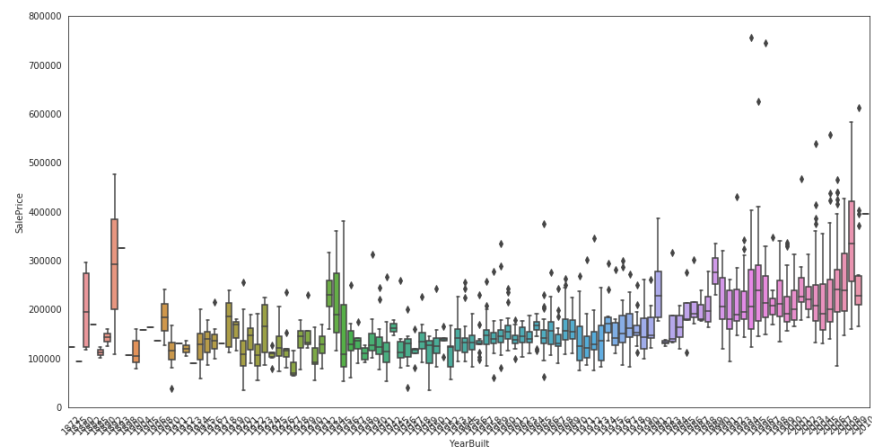
In [9]:

```
data = pd.concat([train['SalePrice'], train['OverallQual']], axis=
1)
f, ax = plt.subplots(figsize=(8, 6))
fig = sns.boxplot(x=train['OverallQual'], y="SalePrice", data=data
)
fig.axis(ymin=0, ymax=800000);
```



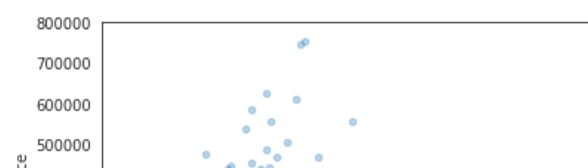
In [10]:

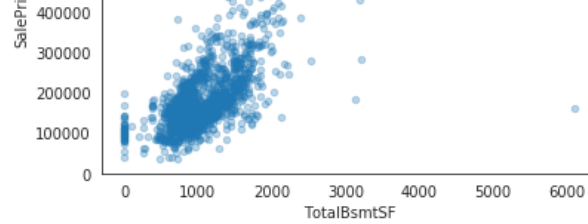
```
data = pd.concat([train['SalePrice'], train['YearBuilt']], axis=1)
f, ax = plt.subplots(figsize=(16, 8))
fig = sns.boxplot(x=train['YearBuilt'], y="SalePrice", data=data)
fig.axis(ymin=0, ymax=800000);
plt.xticks(rotation=45);
```



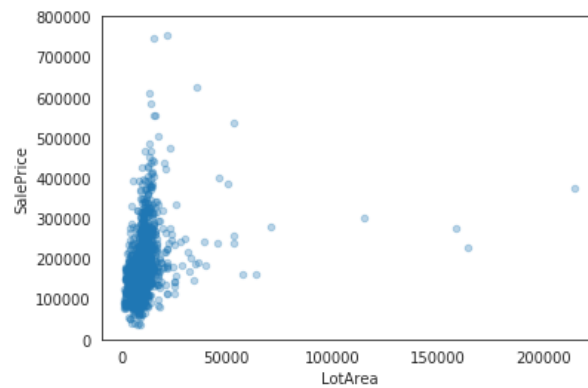
In [11]:

```
data = pd.concat([train['SalePrice'], train['TotalBsmtSF']], axis=
1)
data.plot.scatter(x='TotalBsmtSF', y='SalePrice', alpha=0.3, ylim=
(0,800000));
```

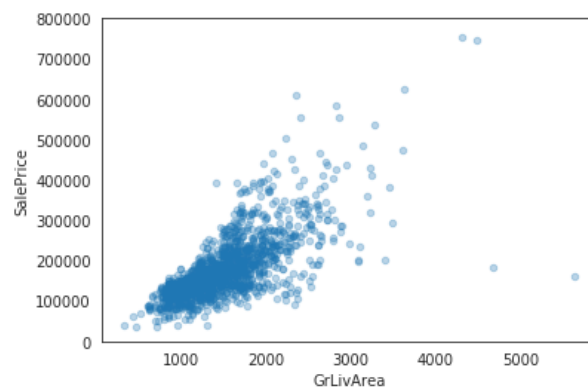




```
In [12]: data = pd.concat([train['SalePrice'], train['LotArea']], axis=1)
data.plot.scatter(x='LotArea', y='SalePrice', alpha=0.3, ylim=(0,800000));
```



```
In [13]: data = pd.concat([train['SalePrice'], train['GrLivArea']], axis=1)
data.plot.scatter(x='GrLivArea', y='SalePrice', alpha=0.3, ylim=(0,800000));
```



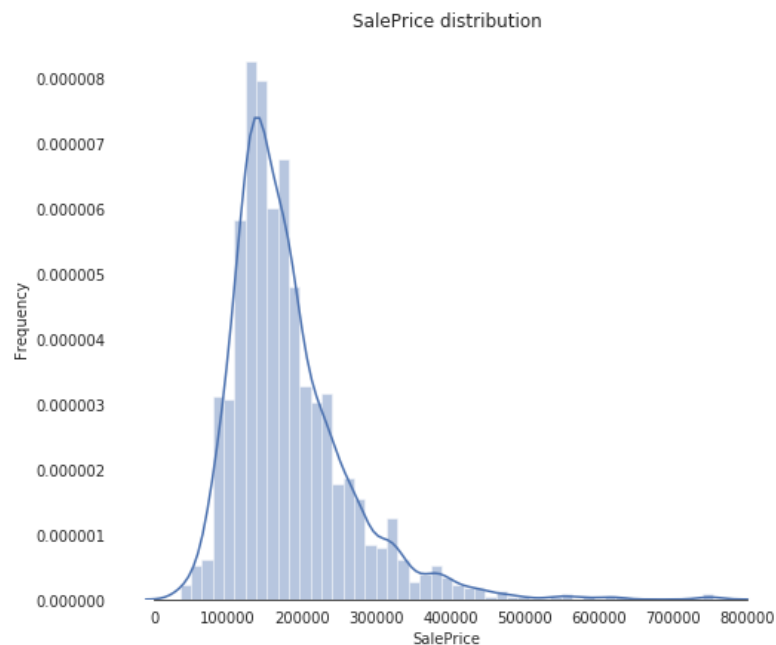
```
In [14]: # Remove the Ids from train and test, as they are unique for each row and hence not useful for the model
train_ID = train['Id']
test_ID = test['Id']
train.drop(['Id'], axis=1, inplace=True)
test.drop(['Id'], axis=1, inplace=True)
train.shape, test.shape
```

```
Out[14]: ((1460, 80), (1459, 79))
```



Let's take a look at the distribution of the SalePrice.

```
In [15]: sns.set_style("white")
sns.set_color_codes(palette='deep')
f, ax = plt.subplots(figsize=(8, 7))
#Check the new distribution
sns.distplot(train['SalePrice'], color="b");
ax.xaxis.grid(False)
ax.set(ylabel="Frequency")
ax.set(xlabel="SalePrice")
ax.set(title="SalePrice distribution")
sns.despine(trim=True, left=True)
plt.show()
```



The SalePrice is skewed to the right. This is a problem because most ML models don't do well with non-normally distributed data. We can apply a  $\log(1+x)$  transform to fix the skew.

```
In [16]: # log(1+x) transform
train["SalePrice"] = np.log1p(train["SalePrice"])
```

Let's plot the SalePrice again.

```
In [17]: sns.set_style("white")
sns.set_color_codes(palette='deep')
f, ax = plt.subplots(figsize=(8, 7))
#Check the new distribution
sns.distplot(train['SalePrice'], fit=norm, color="b");
```

```

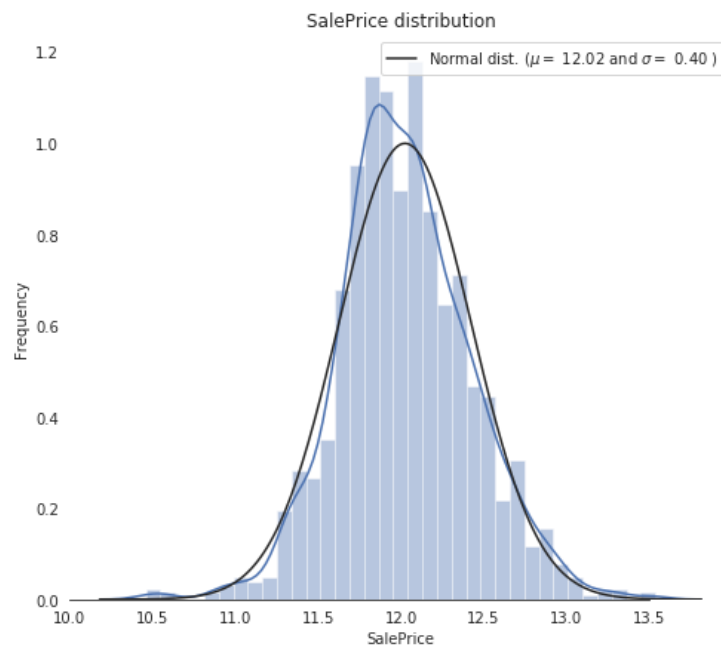
# Get the fitted parameters used by the function
(mu, sigma) = norm.fit(train['SalePrice'])
print( '\n mu = {:.2f} and sigma = {:.2f}\n'.format(mu, sigma))

#Now plot the distribution
plt.legend(['Normal dist. ( $\mu$ = $\$ {:.2f}$  and  $\sigma$ = $\$ {:.2f}$ )'].format(mu, sigma)),
          loc='best')
ax.xaxis.grid(False)
ax.set(ylabel="Frequency")
ax.set(xlabel="SalePrice")
ax.set(title="SalePrice distribution")
sns.despine(trim=True, left=True)

plt.show()

```

mu = 12.02 and sigma = 0.40



The SalePrice is now normally distributed, excellent!

```

In [18]: # Remove outliers
train.drop(train[(train['OverallQual']<5) & (train['SalePrice']>20
0000)].index, inplace=True)
train.drop(train[(train['GrLivArea']>4500) & (train['SalePrice']<3
00000)].index, inplace=True)
train.reset_index(drop=True, inplace=True)

```

```

In [19]: # Split features and labels
train_labels = train['SalePrice'].reset_index(drop=True)
train_features = train.drop(['SalePrice'], axis=1)
test_features = test

# Combine train and test features in order to apply the feature tran
sformation pipeline to the entire dataset

```

```
all_features = pd.concat([train_features, test_features]).reset_index(drop=True)
all_features.shape
```

```
Out[19]:
(2917, 79)
```

## Fill missing values

```
In [20]:
# determine the threshold for missing values
def percent_missing(df):
    data = pd.DataFrame(df)
    df_cols = list(pd.DataFrame(data))
    dict_x = {}
    for i in range(0, len(df_cols)):
        dict_x.update({df_cols[i]: round(data[df_cols[i]].isnull().mean()*100,2)})

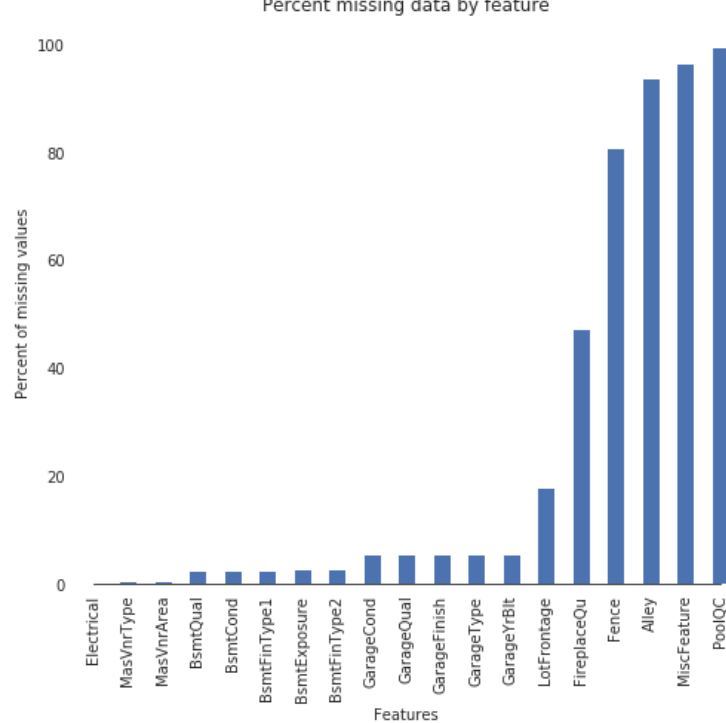
    return dict_x

missing = percent_missing(all_features)
df_miss = sorted(missing.items(), key=lambda x: x[1], reverse=True)
print('Percent of missing data')
df_miss[0:10]
```

Percent of missing data

```
Out[20]:
[('PoolQC', 99.69),
 ('MiscFeature', 96.4),
 ('Alley', 93.21),
 ('Fence', 80.43),
 ('FireplaceQu', 48.68),
 ('LotFrontage', 16.66),
 ('GarageYrBlt', 5.45),
 ('GarageFinish', 5.45),
 ('GarageQual', 5.45),
 ('GarageCond', 5.45)]
```

```
In [21]:
# Visualize missing values
sns.set_style("white")
f, ax = plt.subplots(figsize=(8, 7))
sns.set_color_codes(palette='deep')
missing = round(train.isnull().mean()*100,2)
missing = missing[missing > 0]
missing.sort_values(inplace=True)
missing.plot.bar(color="b")
# Tweak the visual presentation
ax.xaxis.grid(False)
ax.set(ylabel="Percent of missing values")
ax.set(xlabel="Features")
ax.set(title="Percent missing data by feature")
sns.despine(trim=True, left=True)
```



We can now move through each of the features above and impute the missing values for each of them.

```
In [22]: # Some of the non-numeric predictors are stored as numbers; convert them into strings
all_features['MSSubClass'] = all_features['MSSubClass'].apply(str)
all_features['YrSold'] = all_features['YrSold'].astype(str)
all_features['MoSold'] = all_features['MoSold'].astype(str)
```

```
In [23]: def handle_missing(features):
    # the data description states that NA refers to typical ('Typ') values
    features['Functional'] = features['Functional'].fillna('Typ')
    # Replace the missing values in each of the columns below with their mode
    features['Electrical'] = features['Electrical'].fillna("SBrkr")

    features['KitchenQual'] = features['KitchenQual'].fillna("TA")
    features['Exterior1st'] = features['Exterior1st'].fillna(features['Exterior1st'].mode()[0])
    features['Exterior2nd'] = features['Exterior2nd'].fillna(features['Exterior2nd'].mode()[0])
    features['SaleType'] = features['SaleType'].fillna(features['SaleType'].mode()[0])
    features['MSZoning'] = features.groupby('MSSubClass')['MSZoning'].transform(lambda x: x.fillna(x.mode()[0]))

    # the data description stats that NA refers to "No Pool"
    features["PoolQC"] = features["PoolQC"].fillna("None")
    # Replacing the missing values with 0, since no garage = no cars in garage
    for col in ('GarageYrBlt', 'GarageArea', 'GarageCars'):
        features[col] = features[col].fillna(0)
    # Replacing the missing values with None
    for col in ['GarageType', 'GarageFinish', 'GarageQual', 'GarageYrBlt', 'GarageArea', 'GarageCars', 'GarageType', 'GarageFinish', 'GarageQual', 'GarageYrBlt', 'GarageArea', 'GarageCars']:
```

```

eCond']:
    features[col] = features[col].fillna('None')
    # NaN values for these categorical basement features, means there's no basement
    for col in ('BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2'):
        features[col] = features[col].fillna('None')

    # Group the by neighborhoods, and fill in missing value by the median LotFrontage of the neighborhood
    features['LotFrontage'] = features.groupby('Neighborhood')['LotFrontage'].transform(lambda x: x.fillna(x.median()))

    # We have no particular intuition around how to fill in the rest of the categorical features
    # So we replace their missing values with None
    objects = []
    for i in features.columns:
        if features[i].dtype == object:
            objects.append(i)
    features.update(features[objects].fillna('None'))

    # And we do the same thing for numerical features, but this time with 0s
    numeric_dtypes = ['int16', 'int32', 'int64', 'float16', 'float32', 'float64']
    numeric = []
    for i in features.columns:
        if features[i].dtype in numeric_dtypes:
            numeric.append(i)
    features.update(features[numeric].fillna(0))
    return features

all_features = handle_missing(all_features)

```

```

In [24]: # Let's make sure we handled all the missing values
missing = percent_missing(all_features)
df_miss = sorted(missing.items(), key=lambda x: x[1], reverse=True)
print('Percent of missing data')
df_miss[0:10]

```

Percent of missing data

```

Out[24]:
[('MSSubClass', 0.0),
 ('MSZoning', 0.0),
 ('LotFrontage', 0.0),
 ('LotArea', 0.0),
 ('Street', 0.0),
 ('Alley', 0.0),
 ('LotShape', 0.0),
 ('LandContour', 0.0),
 ('Utilities', 0.0),
 ('LotConfig', 0.0)]

```

There are no missing values anymore!

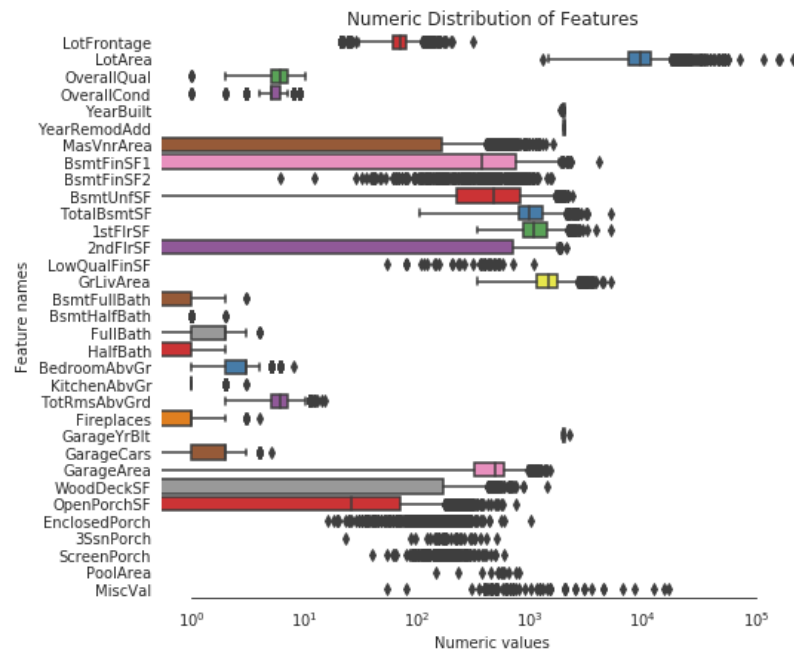
## Fix skewed features

In [25]:

```
# Fetch all numeric features
numeric_dtypes = ['int16', 'int32', 'int64', 'float16', 'float32',
                  'float64']
numeric = []
for i in all_features.columns:
    if all_features[i].dtype in numeric_dtypes:
        numeric.append(i)
```

In [26]:

```
# Create box plots for all numeric features
sns.set_style("white")
f, ax = plt.subplots(figsize=(8, 7))
ax.set_xscale("log")
ax = sns.boxplot(data=all_features[numeric], orient="h", palette="Set1")
ax.xaxis.grid(False)
ax.set(ylabel="Feature names")
ax.set(xlabel="Numeric values")
ax.set(title="Numeric Distribution of Features")
sns.despine(trim=True, left=True)
```



In [27]:

```
# Find skewed numerical features
skew_features = all_features[numeric].apply(lambda x: skew(x)).sort_values(ascending=False)

high_skew = skew_features[skew_features > 0.5]
skew_index = high_skew.index

print("There are {} numerical features with Skew > 0.5 :".format(high_skew.shape[0]))
skewness = pd.DataFrame({'Skew' :high_skew})
skew_features.head(10)
```

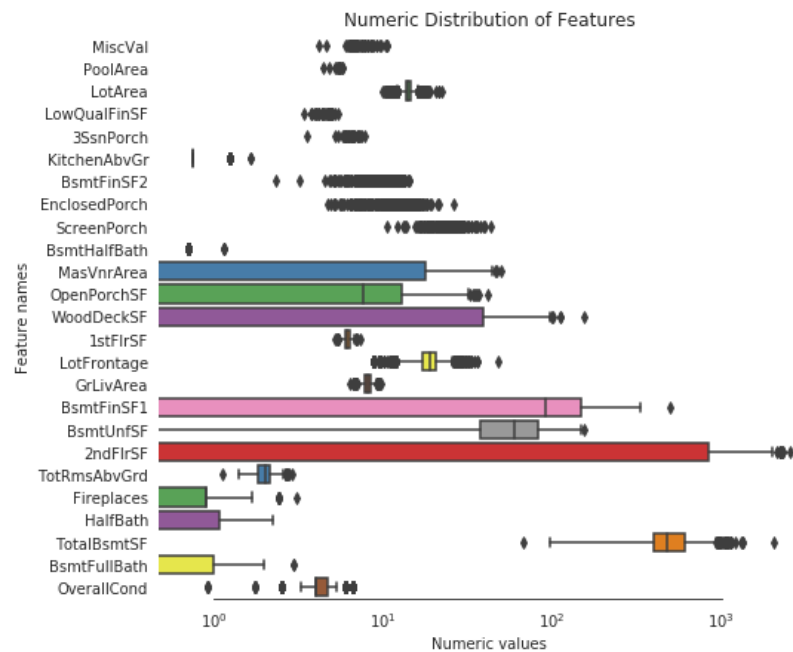
There are 25 numerical features with Skew > 0.5 :

```
Out[27]:
MiscVal      21.939672
PoolArea     17.688664
LotArea      13.109495
LowQualFinSF 12.084539
3SsnPorch    11.372080
KitchenAbvGr 4.300550
BsmtFinSF2   4.144503
EnclosedPorch 4.002344
ScreenPorch   3.945101
BsmtHalfBath  3.929996
dtype: float64
```

We use the scipy function `boxcox1p` which computes the Box-Cox transformation. The goal is to find a simple transformation that lets us normalize data.

```
In [28]:
# Normalize skewed features
for i in skew_index:
    all_features[i] = boxcox1p(all_features[i], boxcox_normmax(all_
    _features[i] + 1))
```

```
In [29]:
# Let's make sure we handled all the skewed values
sns.set_style("white")
f, ax = plt.subplots(figsize=(8, 7))
ax.set_xscale("log")
ax = sns.boxplot(data=all_features[skew_index], orient="h", palette="Set1")
ax.xaxis.grid(False)
ax.set(ylabel="Feature names")
ax.set(xlabel="Numeric values")
ax.set(title="Numeric Distribution of Features")
sns.despine(trim=True, left=True)
```



All the features look fairly normally distributed now.

## Create interesting features

ML models have trouble recognizing more complex patterns (and we're staying away from neural nets for this competition), so let's help our models out by creating a few features based on our intuition about the dataset, e.g. total area of floors, bathrooms and porch area of each house.

```
In [30]: all_features['BsmtFinType1_Unf'] = 1*(all_features['BsmtFinType1']
== 'Unf')
all_features['HasWoodDeck'] = (all_features['WoodDeckSF'] == 0) *
1
all_features['HasOpenPorch'] = (all_features['OpenPorchSF'] == 0)
* 1
all_features['HasEnclosedPorch'] = (all_features['EnclosedPorch']
== 0) * 1
all_features['Has3SsnPorch'] = (all_features['3SsnPorch'] == 0) *
1
all_features['HasScreenPorch'] = (all_features['ScreenPorch'] == 0
) * 1
all_features['YearsSinceRemodel'] = all_features['YrSold'].astype(
int) - all_features['YearRemodAdd'].astype(int)
all_features['Total_Home_Quality'] = all_features['OverallQual'] +
all_features['OverallCond']
all_features = all_features.drop(['Utilities', 'Street', 'PoolQC'
, ], axis=1)
all_features['TotalSF'] = all_features['TotalBsmtSF'] + all_featur
es['1stFlrSF'] + all_features['2ndFlrSF']
all_features['YrBltAndRemod'] = all_features['YearBuilt'] + all_fe
atures['YearRemodAdd']

all_features['Total_sqr_footage'] = (all_features['BsmtFinSF1'] +
all_features['BsmtFinSF2'] +
                                all_features['1stFlrSF'] + all_fe
atures['2ndFlrSF'])
all_features['Total_Bathrooms'] = (all_features['FullBath'] + (0.5
* all_features['HalfBath']) +
                                all_features['BsmtFullBath'] + (0.5
* all_features['BsmtHalfBath']))
all_features['Total_porch_sf'] = (all_features['OpenPorchSF'] + al
l_features['3SsnPorch'] +
                                all_features['EnclosedPorch'] + all_
features['ScreenPorch'] +
                                all_features['WoodDeckSF'])
all_features['TotalBsmtSF'] = all_features['TotalBsmtSF'].apply(la
mbda x: np.exp(6) if x <= 0.0 else x)
all_features['2ndFlrSF'] = all_features['2ndFlrSF'].apply(lambda x
: np.exp(6.5) if x <= 0.0 else x)
all_features['GarageArea'] = all_features['GarageArea'].apply(lamb
da x: np.exp(6) if x <= 0.0 else x)
all_features['GarageCars'] = all_features['GarageCars'].apply(lamb
da x: 0 if x <= 0.0 else x)
```



```

all_features['LotFrontage'] = all_features['LotFrontage'].apply(lambda x: np.exp(4.2) if x <= 0.0 else x)
all_features['MasVnrArea'] = all_features['MasVnrArea'].apply(lambda x: np.exp(4) if x <= 0.0 else x)
all_features['BsmtFinSF1'] = all_features['BsmtFinSF1'].apply(lambda x: np.exp(6.5) if x <= 0.0 else x)

all_features['haspool'] = all_features['PoolArea'].apply(lambda x: 1 if x > 0 else 0)
all_features['has2ndfloor'] = all_features['2ndFlrSF'].apply(lambda x: 1 if x > 0 else 0)
all_features['hasgarage'] = all_features['GarageArea'].apply(lambda x: 1 if x > 0 else 0)
all_features['hasbsmt'] = all_features['TotalBsmtSF'].apply(lambda x: 1 if x > 0 else 0)
all_features['hasfireplace'] = all_features['Fireplaces'].apply(lambda x: 1 if x > 0 else 0)

```

## Feature transformations

Let's create more features by calculating the log and square transformations of our numerical features. We do this manually, because ML models won't be able to reliably tell if  $\log(\text{feature})$  or  $\text{feature}^2$  is a predictor of the SalePrice.

```

In [31]:
def logs(res, ls):
    m = res.shape[1]
    for l in ls:
        res = res.assign(newcol=pd.Series(np.log(1.01+res[l])).values)
        res.columns.values[m] = l + '_log'
        m += 1
    return res

log_features = ['LotFrontage', 'LotArea', 'MasVnrArea', 'BsmtFinSF1',
                'BsmtFinSF2', 'BsmtUnfSF',
                'TotalBsmtSF', '1stFlrSF', '2ndFlrSF', 'LowQualFinSF',
                'GrLivArea',
                'BsmtFullBath', 'BsmtHalfBath', 'FullBath', 'HalfBath',
                'BedroomAbvGr', 'KitchenAbvGr',
                'TotRmsAbvGrd', 'Fireplaces', 'GarageCars', 'GarageArea',
                'WoodDeckSF', 'OpenPorchSF',
                'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'PoolArea',
                'MiscVal', 'YearRemodAdd', 'TotalSF']

all_features = logs(all_features, log_features)

```

```

In [32]:
def squares(res, ls):
    m = res.shape[1]
    for l in ls:
        res = res.assign(newcol=pd.Series(res[l]*res[l]).values)
        res.columns.values[m] = l + '_sq'
        m += 1
    return res

squared_features = ['YearRemodAdd', 'LotFrontage_log',

```

```

    'TotalBsmtSF_log', '1stFlrSF_log', '2ndFlrSF_log',
    'GrLivArea_log',
    'GarageCars_log', 'GarageArea_log']
all_features = squares(all_features, squared_features)

```

## Encode categorical features

Numerically encode categorical features because most models can only handle numerical features.

```

In [33]: all_features = pd.get_dummies(all_features).reset_index(drop=True)
all_features.shape

```

```

Out[33]:
(2917, 379)

```

```

In [34]: all_features.head()

```

```

Out[34]:

```

|   | LotFrontage | LotArea   | OverallQual | OverallCond | YearBuilt | YearRemodAdd | MasVnrArea |
|---|-------------|-----------|-------------|-------------|-----------|--------------|------------|
| 0 | 18.144572   | 13.833055 | 7           | 3.991517    | 2003      | 2003         | 19.433174  |
| 1 | 20.673625   | 14.117918 | 6           | 6.000033    | 1976      | 1976         | 54.598150  |
| 2 | 18.668046   | 14.476513 | 7           | 3.991517    | 2001      | 2002         | 17.768840  |
| 3 | 17.249650   | 14.106197 | 7           | 3.991517    | 1915      | 1970         | 54.598150  |
| 4 | 21.314282   | 15.022009 | 8           | 3.991517    | 2000      | 2000         | 25.404163  |

```

In [35]: all_features.shape

```

```

Out[35]:
(2917, 379)

```

```

In [36]: # Remove any duplicated column names
all_features = all_features.loc[:,~all_features.columns.duplicated
()]

```

## Recreate training and test sets

```

In [37]: X = all_features.iloc[:len(train_labels), :]
X_test = all_features.iloc[len(train_labels):, :]
X.shape, train_labels.shape, X_test.shape

```

Out[37]:

```
((1458, 378), (1458,), (1459, 378))
```

Visualize some of the features we're going to train our models on.

In [38]:

```
# Finding numeric features
numeric_dtypes = ['int16', 'int32', 'int64', 'float16', 'float32',
                  'float64']
numeric = []
for i in X.columns:
    if X[i].dtype in numeric_dtypes:
        if i in ['TotalSF', 'Total_Bathrooms', 'Total_porch_sf', 'haspool', 'hasgarage', 'hasbsmt', 'hasfireplace']:
            pass
        else:
            numeric.append(i)

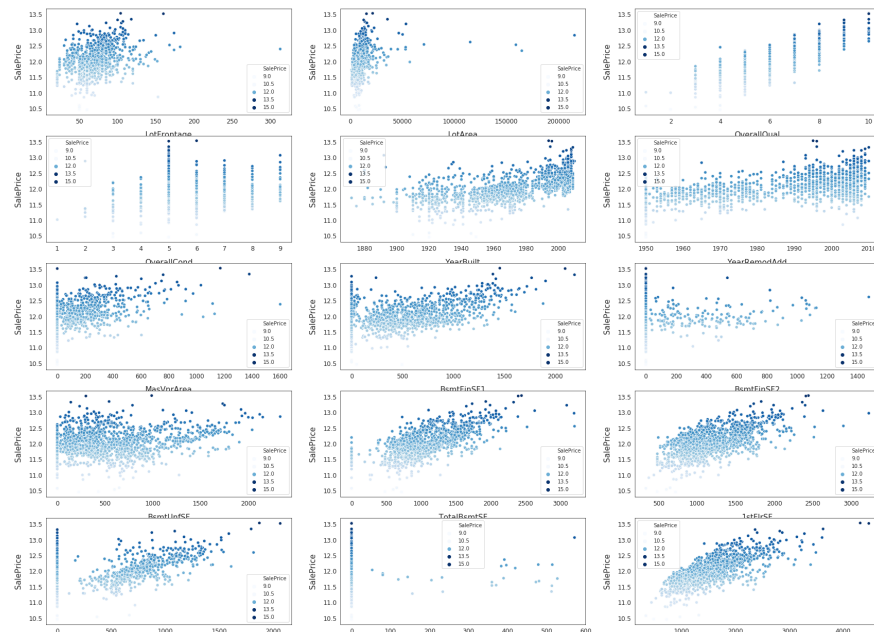
# visualising some more outliers in the data values
fig, axs = plt.subplots(ncols=2, nrows=0, figsize=(12, 150))
plt.subplots_adjust(right=2)
plt.subplots_adjust(top=2)
sns.color_palette("husl", 8)
for i, feature in enumerate(list(X[numeric]), 1):
    if(feature=='MiscVal'):
        break
    plt.subplot(len(list(numeric)), 3, i)
    sns.scatterplot(x=feature, y='SalePrice', hue='SalePrice', palette='Blues', data=train)

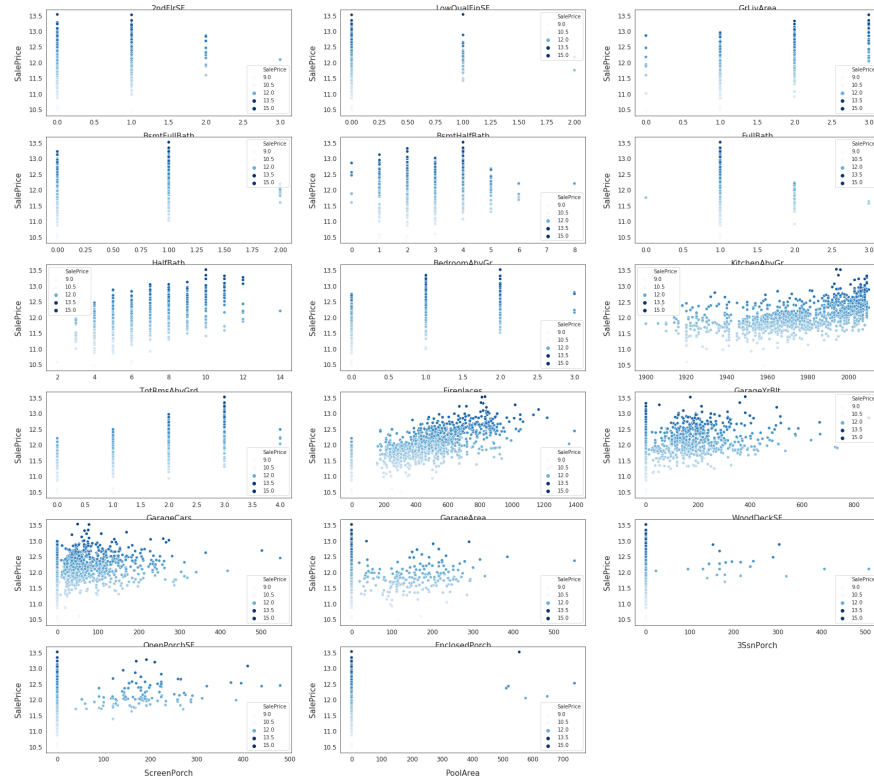
    plt.xlabel('{}'.format(feature), size=15, labelpad=12.5)
    plt.ylabel('SalePrice', size=15, labelpad=12.5)

    for j in range(2):
        plt.tick_params(axis='x', labelsize=12)
        plt.tick_params(axis='y', labelsize=12)

    plt.legend(loc='best', prop={'size': 10})

plt.show()
```





## Train a model

### Key features of the model training process:

- **Cross Validation:** Using 12-fold cross-validation
- **Models:** On each run of cross-validation I fit 7 models (ridge, svr, gradient boosting, random forest, xgboost, lightgbm regressors)
- **Stacking:** In addition, I trained a meta StackingCVRegressor optimized using xgboost
- **Blending:** All models trained will overfit the training data to varying degrees. Therefore, to make final predictions, I blended their predictions together to get more robust predictions.

### Setup cross validation and define error metrics

```
In [39]: # Setup cross validation folds
kf = KFold(n_splits=12, random_state=42, shuffle=True)
```

```
In [40]: # Define error metrics
def rmsle(y, y_pred):
    return np.sqrt(mean_squared_error(y, y_pred))

def cv_rmse(model, X=X):
```

```

        rmse = np.sqrt(-cross_val_score(model, X, train_labels, scoring="neg_mean_squared_error", cv=kf))
    return (rmse)

```

## Setup models

In [41]:

```

# Light Gradient Boosting Regressor
lightgbm = LGBMRegressor(objective='regression',
                           num_leaves=6,
                           learning_rate=0.01,
                           n_estimators=7000,
                           max_bin=200,
                           bagging_fraction=0.8,
                           bagging_freq=4,
                           bagging_seed=8,
                           feature_fraction=0.2,
                           feature_fraction_seed=8,
                           min_sum_hessian_in_leaf = 11,
                           verbose=-1,
                           random_state=42)

# XGBoost Regressor
xgboost = XGBRegressor(learning_rate=0.01,
                        n_estimators=6000,
                        max_depth=4,
                        min_child_weight=0,
                        gamma=0.6,
                        subsample=0.7,
                        colsample_bytree=0.7,
                        objective='reg:linear',
                        nthread=-1,
                        scale_pos_weight=1,
                        seed=27,
                        reg_alpha=0.00006,
                        random_state=42)

# Ridge Regressor
ridge_alphas = [1e-15, 1e-10, 1e-8, 9e-4, 7e-4, 5e-4, 3e-4, 1e-4,
                1e-3, 5e-2, 1e-2, 0.1, 0.3, 1, 3, 5, 10, 15, 18, 20, 30, 50, 75, 100]
ridge = make_pipeline(RobustScaler(), RidgeCV(alphas=ridge_alphas,
                                              cv=kf))

# Support Vector Regressor
svr = make_pipeline(RobustScaler(), SVR(C= 20, epsilon= 0.008, gamma=0.0003))

# Gradient Boosting Regressor
gbr = GradientBoostingRegressor(n_estimators=6000,
                                 learning_rate=0.01,
                                 max_depth=4,
                                 max_features='sqrt',
                                 min_samples_leaf=15,
                                 min_samples_split=10,
                                 loss='huber',
                                 random_state=42)

```

```
# Random Forest Regressor
rf = RandomForestRegressor(n_estimators=1200,
                           max_depth=15,
                           min_samples_split=5,
                           min_samples_leaf=5,
                           max_features=None,
                           oob_score=True,
                           random_state=42)

# Stack up all the models above, optimized using xgboost
stack_gen = StackingCVRegressor(regressors=(xgboost, lightgbm, svr
, ridge, gbr, rf),
                                meta_regressor=xgboost,
                                use_features_in_secondary=True)
```

## Train models

Get cross validation scores for each model

```
In [42]:
scores = {}

score = cv_rmse(lightgbm)
print("lightgbm: {:.4f} ({:.4f})".format(score.mean(), score.std
()))
scores['lgb'] = (score.mean(), score.std())

lightgbm: 0.1159 (0.0167)
```

```
In [43]:
score = cv_rmse(xgboost)
print("xgboost: {:.4f} ({:.4f})".format(score.mean(), score.std
()))
scores['xgb'] = (score.mean(), score.std())

xgboost: 0.1364 (0.0175)
```

```
In [44]:
score = cv_rmse(svr)
print("SVR: {:.4f} ({:.4f})".format(score.mean(), score.std()))
scores['svr'] = (score.mean(), score.std())

SVR: 0.1094 (0.0200)
```

```
In [45]:
score = cv_rmse(ridge)
print("ridge: {:.4f} ({:.4f})".format(score.mean(), score.std()))
scores['ridge'] = (score.mean(), score.std())

ridge: 0.1101 (0.0161)
```

```
In [46]: score = cv_rmse(rf)
print("rf: {:.4f} ({:.4f})".format(score.mean(), score.std()))
scores['rf'] = (score.mean(), score.std())

rf: 0.1366 (0.0188)
```



## How I made top 0.3% on a Kaggle competition

Python notebook using data from [multiple data sources](#) · 18,762 views · 3d ago · starter code, eda, feature engineering

, +3 more



88

Copy and Edit

94

### Version 31

[31 commits](#)

### Notebook

How I Made Top 0.3%  
On A Kaggle  
Competition

EDA

Feature Engineering

[Train A Model](#)

Data

Output

Log

Comments

gbr: 0.1121 (0.0164)

### Fit the models

```
In [48]: print('stack_gen')
stack_gen_model = stack_gen.fit(np.array(X), np.array(train_labels))
```

stack\_gen

```
In [49]: print('lightgbm')
lgb_model_full_data = lightgbm.fit(X, train_labels)
```

lightgbm

```
In [50]: print('xgboost')
xgb_model_full_data = xgboost.fit(X, train_labels)
```

xgboost

```
In [51]: print('Svr')
svr_model_full_data = svr.fit(X, train_labels)
```

Svr

```
In [52]: print('Ridge')
ridge_model_full_data = ridge.fit(X, train_labels)
```

Ridge

```
In [53]: print('RandomForest')
rf_model_full_data = rf.fit(X, train_labels)
```

```
In [54]: print('GradientBoosting')
gbr_model_full_data = gbr.fit(X, train_labels)
```



## Blend models and get predictions

```
In [55]: # Blend models in order to make the final predictions more robust to overfitting
def blended_predictions(X):
    return ((0.1 * ridge_model_full_data.predict(X)) + \
            (0.2 * svr_model_full_data.predict(X)) + \
            (0.1 * gbr_model_full_data.predict(X)) + \
            (0.1 * xgb_model_full_data.predict(X)) + \
            (0.1 * lgb_model_full_data.predict(X)) + \
            (0.05 * rf_model_full_data.predict(X)) + \
            (0.35 * stack_gen_model.predict(np.array(X))))
```

```
In [56]: # Get final predictions from the blended model
blended_score = rmsle(train_labels, blended_predictions(X))
scores['blended'] = (blended_score, 0)
print('RMSLE score on train data:')
print(blended_score)
```

```
RMSLE score on train data:
0.07537440195302639
```

## Identify the best performing model

```
In [57]: # Plot the predictions for each model
sns.set_style("white")
fig = plt.figure(figsize=(24, 12))

ax = sns.pointplot(x=list(scores.keys()), y=[score for score, _ in
scores.values()], markers=['o'], linestyle=['-'])
for i, score in enumerate(scores.values()):
    ax.text(i, score[0] + 0.002, '{:.6f}'.format(score[0]), horizontalalignment='left', size='large', color='black', weight='semibold')

plt.ylabel('Score (RMSE)', size=20, labelpad=12.5)
```