

# COMPREHENSIVE DATA EXPLORATION WITH PYTHON

Pedro Marcelino (<http://pmarcelino.com>) - February 2017

Other Kernels: [Data analysis and feature extraction with Python \(https://www.kaggle.com/pmarcelino/data-analysis-and-feature-extraction-with-python\)](https://www.kaggle.com/pmarcelino/data-analysis-and-feature-extraction-with-python)

---

## 'The most difficult thing in life is to know yourself'

This quote belongs to Thales of Miletus. Thales was a Greek/Phonician philosopher, mathematician and astronomer, which is recognised as the first individual in Western civilisation known to have entertained and engaged in scientific thought (source: <https://en.wikipedia.org/wiki/Thales> (<https://en.wikipedia.org/wiki/Thales>))

I wouldn't say that knowing your data is the most difficult thing in data science, but it is time-consuming. Therefore, it's easy to overlook this initial step and jump too soon into the water.

So I tried to learn how to swim before jumping into the water. Based on [Hair et al. \(2013\) \(https://amzn.to/2JuDmvo\)](https://amzn.to/2JuDmvo), chapter 'Examining your data', I did my best to follow a comprehensive, but not exhaustive, analysis of the data. I'm far from reporting a rigorous study in this kernel, but I hope that it can be useful for the community, so I'm sharing how I applied some of those data analysis principles to this problem.

Despite the strange names I gave to the chapters, what we are doing in this kernel is something like:

1. **Understand the problem.** We'll look at each variable and do a philosophical analysis about their meaning and importance for this problem.
2. **Univariable study.** We'll just focus on the dependent variable ('SalePrice') and try to know a little bit more about it.
3. **Multivariate study.** We'll try to understand how the dependent variable and independent variables relate.
4. **Basic cleaning.** We'll clean the dataset and handle the missing data, outliers and categorical variables.
5. **Test assumptions.** We'll check if our data meets the assumptions required by most multivariate techniques.

Now, it's time to have fun!

```
In [1]: #invite people for the Kaggle party
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from scipy.stats import norm
from sklearn.preprocessing import StandardScaler
from scipy import stats
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

```
In [2]: #bring in the six packs
df_train = pd.read_csv('../input/train.csv')
```

```
In [3]: #check the decoration  
df_train.columns
```

```
Out[3]: Index(['Id', 'MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street',  
'Alley', 'LotShape', 'LandContour', 'Utilities', 'LotConfig',  
'LandSlope', 'Neighborhood', 'Condition1', 'Condition2', 'BldgType',  
'HouseStyle', 'OverallQual', 'OverallCond', 'YearBuilt', 'YearRemodAdd',  
,  
'RoofStyle', 'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType',  
'MasVnrArea', 'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual',  
'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinSF1',  
'BsmtFinType2', 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', 'Heating',  
'HeatingQC', 'CentralAir', 'Electrical', '1stFlrSF', '2ndFlrSF',  
'LowQualFinSF', 'GrLivArea', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath',  
,  
'HalfBath', 'BedroomAbvGr', 'KitchenAbvGr', 'KitchenQual',  
'TotRmsAbvGrd', 'Functional', 'Fireplaces', 'FireplaceQu', 'GarageType',  
,  
'GarageYrBlt', 'GarageFinish', 'GarageCars', 'GarageArea', 'GarageQual',  
,  
'GarageCond', 'PavedDrive', 'WoodDeckSF', 'OpenPorchSF',  
'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'PoolArea', 'PoolQC',  
'Fence', 'MiscFeature', 'MiscVal', 'MoSold', 'YrSold', 'SaleType',  
'SaleCondition', 'SalePrice'],  
dtype='object')
```

# 1. So... What can we expect?

In order to understand our data, we can look at each variable and try to understand their meaning and relevance to this problem. I know this is time-consuming, but it will give us the flavour of our dataset.

In order to have some discipline in our analysis, we can create an Excel spreadsheet with the following columns:

- **Variable** - Variable name.
- **Type** - Identification of the variables' type. There are two possible values for this field: 'numerical' or 'categorical'. By 'numerical' we mean variables for which the values are numbers, and by 'categorical' we mean variables for which the values are categories.
- **Segment** - Identification of the variables' segment. We can define three possible segments: building, space or location. When we say 'building', we mean a variable that relates to the physical characteristics of the building (e.g. 'OverallQual'). When we say 'space', we mean a variable that reports space properties of the house (e.g. 'TotalBsmstSF'). Finally, when we say a 'location', we mean a variable that gives information about the place where the house is located (e.g. 'Neighborhood').
- **Expectation** - Our expectation about the variable influence in 'SalePrice'. We can use a categorical scale with 'High', 'Medium' and 'Low' as possible values.
- **Conclusion** - Our conclusions about the importance of the variable, after we give a quick look at the data. We can keep with the same categorical scale as in 'Expectation'.
- **Comments** - Any general comments that occurred to us.

While 'Type' and 'Segment' is just for possible future reference, the column 'Expectation' is important because it will help us develop a 'sixth sense'. To fill this column, we should read the description of all the variables and, one by one, ask ourselves:

- Do we think about this variable when we are buying a house? (e.g. When we think about the house of our dreams, do we care about its 'Masonry veneer type'?).
- If so, how important would this variable be? (e.g. What is the impact of having 'Excellent' material on the exterior instead of 'Poor'? And of having 'Excellent' instead of 'Good'?).
- Is this information already described in any other variable? (e.g. If 'LandContour' gives the flatness of the property, do we really need to know the 'LandSlope'?).

After this daunting exercise, we can filter the spreadsheet and look carefully to the variables with 'High' 'Expectation'. Then, we can rush into some scatter plots between those variables and 'SalePrice', filling in the 'Conclusion' column which is just the correction of our expectations.

I went through this process and concluded that the following variables can play an important role in this problem:

- OverallQual (which is a variable that I don't like because I don't know how it was computed; a funny exercise would be to predict 'OverallQual' using all the other variables available).
- YearBuilt.
- TotalBsmstSF.
- GrLivArea.

I ended up with two 'building' variables ('OverallQual' and 'YearBuilt') and two 'space' variables ('TotalBsmstSF' and 'GrLivArea'). This might be a little bit unexpected as it goes against the real estate mantra that all that matters is 'location, location and location'. It is possible that this quick data examination process was a bit harsh for categorical variables. For example, I expected the 'Neighborhood' variable to be more relevant, but after the data examination I ended up excluding it. Maybe this is related to the use of scatter plots instead of boxplots, which are more suitable for categorical variables visualization. The way we visualize data often influences our conclusions.

However, the main point of this exercise was to think a little about our data and expectations, so I think we achieved our goal. Now it's time for 'a little less conversation, a little more action please'. Let's **shake it!**

## 2. First things first: analysing 'SalePrice'

'SalePrice' is the reason of our quest. It's like when we're going to a party. We always have a reason to be there. Usually, women are that reason. (disclaimer: adapt it to men, dancing or alcohol, according to your preferences)

Using the women analogy, let's build a little story, the story of 'How we met 'SalePrice'.

*Everything started in our Kaggle party, when we were looking for a dance partner. After a while searching in the dance floor, we saw a girl, near the bar, using dance shoes. That's a sign that she's there to dance. We spend much time doing predictive modelling and participating in analytics competitions, so talking with girls is not one of our super powers. Even so, we gave it a try:*

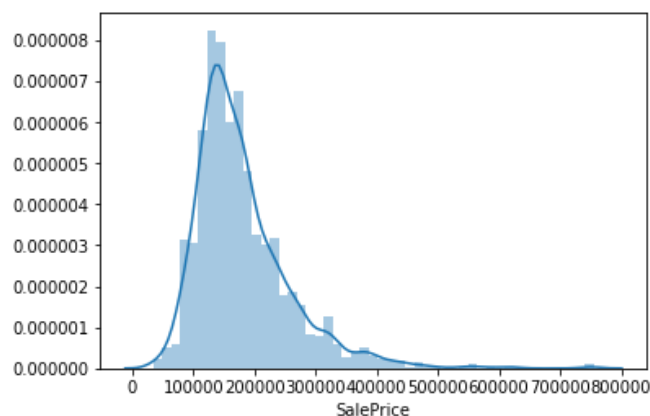
*'Hi, I'm Kaggly! And you? 'SalePrice'? What a beautiful name! You know 'SalePrice', could you give me some data about you? I just developed a model to calculate the probability of a successful relationship between two people. I'd like to apply it to us!'*

```
In [4]: #descriptive statistics summary
df_train['SalePrice'].describe()
```

```
Out[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
```

*'Very well... It seems that your minimum price is larger than zero. Excellent! You don't have one of those personal traits that would destroy my model! Do you have any picture that you can send me? I don't know... like, you in the beach... or maybe a selfie in the gym?'*

```
In [5]: #histogram
sns.distplot(df_train['SalePrice']);
```



'Ah! I see you that you use seaborn makeup when you're going out... That's so elegant! I also see that you:

- **Deviate from the normal distribution.**
- **Have appreciable positive skewness.**
- **Show peakedness.**

This is getting interesting! 'SalePrice', could you give me your body measures?"

```
In [6]: #skewness and kurtosis
print("Skewness: %f" % df_train['SalePrice'].skew())
print("Kurtosis: %f" % df_train['SalePrice'].kurt())

Skewness: 1.882876
Kurtosis: 6.536282
```

'Amazing! If my love calculator is correct, our success probability is 97.834657%. I think we should meet again! Please, keep my number and give me a call if you're free next Friday. See you in a while, crocodile!'

## 'SalePrice', her buddies and her interests

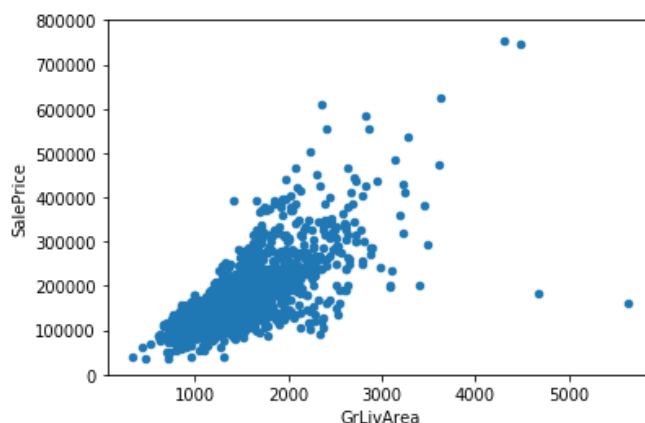
It is military wisdom to choose the terrain where you will fight. As soon as 'SalePrice' walked away, we went to Facebook. Yes, now this is getting serious. Notice that this is not stalking. It's just an intense research of an individual, if you know what I mean.

According to her profile, we have some common friends. Besides Chuck Norris, we both know 'GrLivArea' and 'TotalBsmtSF'. Moreover, we also have common interests such as 'OverallQual' and 'YearBuilt'. This looks promising!

To take the most out of our research, we will start by looking carefully at the profiles of our common friends and later we will focus on our common interests.

### Relationship with numerical variables

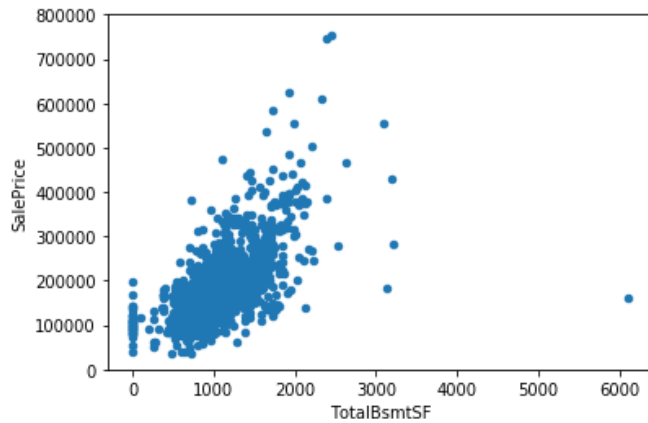
```
In [7]: #scatter plot grlivarea/saleprice
var = 'GrLivArea'
data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
data.plot.scatter(x=var, y='SalePrice', ylim=(0,800000));
```



Hmmm... It seems that 'SalePrice' and 'GrLivArea' are really old friends, with a **linear relationship**.

And what about 'TotalBsmtSF'?

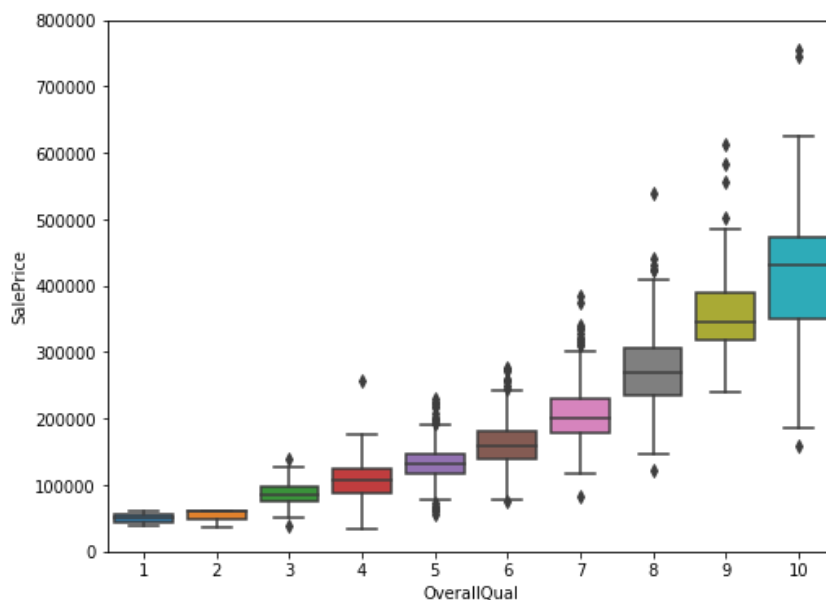
```
In [8]: #scatter plot totalbsmtsf/saleprice
var = 'TotalBsmtSF'
data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
data.plot.scatter(x=var, y='SalePrice', ylim=(0,800000));
```



'TotalBsmtSF' is also a great friend of 'SalePrice' but this seems a much more emotional relationship! Everything is ok and suddenly, in a **strong linear (exponential?)** reaction, everything changes. Moreover, it's clear that sometimes 'TotalBsmtSF' closes in itself and gives zero credit to 'SalePrice'.

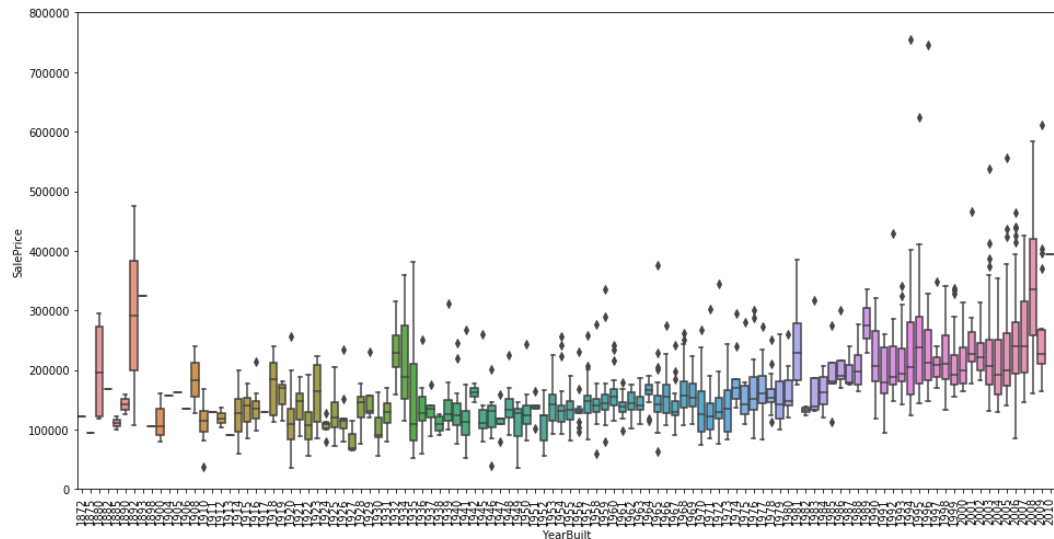
## Relationship with categorical features

```
In [9]: #box plot overallqual/saleprice
var = 'OverallQual'
data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
f, ax = plt.subplots(figsize=(8, 6))
fig = sns.boxplot(x=var, y="SalePrice", data=data)
fig.axis(ymin=0, ymax=800000);
```



*Like all the pretty girls, 'SalePrice' enjoys 'OverallQual'. Note to self: consider whether McDonald's is suitable for the first date.*

```
In [10]: var = 'YearBuilt'
data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
f, ax = plt.subplots(figsize=(16, 8))
fig = sns.boxplot(x=var, y="SalePrice", data=data)
fig.axis(ymin=0, ymax=800000);
plt.xticks(rotation=90);
```



*Although it's not a strong tendency, I'd say that 'SalePrice' is more prone to spend more money in new stuff than in old relics.*

**Note:** we don't know if 'SalePrice' is in constant prices. Constant prices try to remove the effect of inflation. If 'SalePrice' is not in constant prices, it should be, so than prices are comparable over the years.

## In summary

Stories aside, we can conclude that:

- 'GrLivArea' and 'TotalBsmtSF' seem to be linearly related with 'SalePrice'. Both relationships are positive, which means that as one variable increases, the other also increases. In the case of 'TotalBsmtSF', we can see that the slope of the linear relationship is particularly high.
- 'OverallQual' and 'YearBuilt' also seem to be related with 'SalePrice'. The relationship seems to be stronger in the case of 'OverallQual', where the box plot shows how sales prices increase with the overall quality.

We just analysed four variables, but there are many other that we should analyse. The trick here seems to be the choice of the right features (feature selection) and not the definition of complex relationships between them (feature engineering).

That said, let's separate the wheat from the chaff.

## 3. Keep calm and work smart

Until now we just followed our intuition and analysed the variables we thought were important. In spite of our efforts to give an objective character to our analysis, we must say that our starting point was subjective.

As an engineer, I don't feel comfortable with this approach. All my education was about developing a disciplined mind, able to withstand the winds of subjectivity. There's a reason for that. Try to be subjective in structural engineering and you will see physics making things fall down. It can hurt.

So, let's overcome inertia and do a more objective analysis.

## The 'plasma soup'

'In the very beginning there was nothing except for a plasma soup. What is known of these brief moments in time, at the start of our study of cosmology, is largely conjectural. However, science has devised some sketch of what probably happened, based on what is known about the universe today.' (source: <http://umich.edu/~gs265/bigbang.htm>)

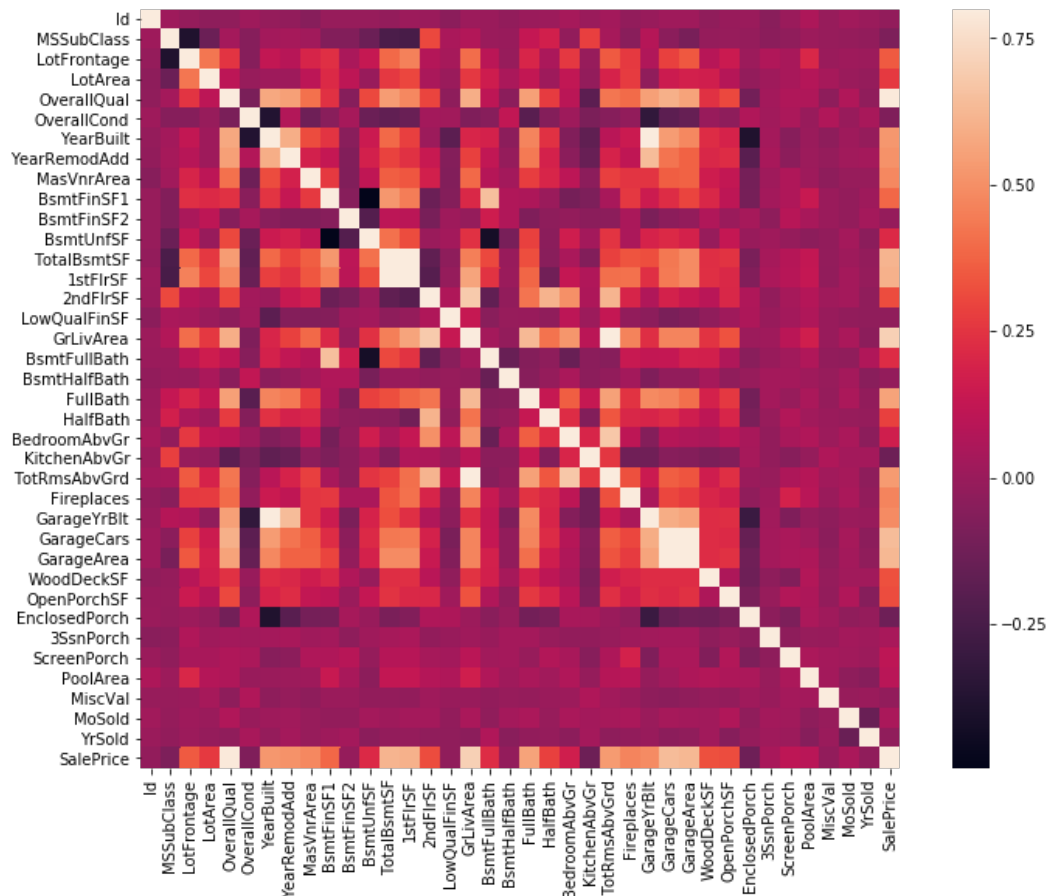
To explore the universe, we will start with some practical recipes to make sense of our 'plasma soup':

- Correlation matrix (heatmap style).
- 'SalePrice' correlation matrix (zoomed heatmap style).
- Scatter plots between the most correlated variables (move like Jagger style).

### Correlation matrix (heatmap style)



```
In [11]: #correlation matrix
corrmat = df_train.corr()
f, ax = plt.subplots(figsize=(12, 9))
sns.heatmap(corrmat, vmax=.8, square=True);
```



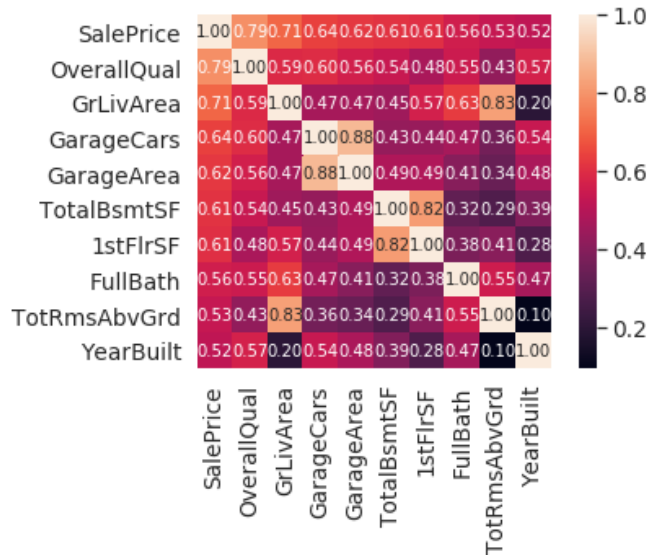
In my opinion, this heatmap is the best way to get a quick overview of our 'plasma soup' and its relationships. (Thank you @seaborn!)

At first sight, there are two red colored squares that get my attention. The first one refers to the 'TotalBsmtSF' and '1stFlrSF' variables, and the second one refers to the 'GarageX' variables. Both cases show how significant the correlation is between these variables. Actually, this correlation is so strong that it can indicate a situation of multicollinearity. If we think about these variables, we can conclude that they give almost the same information so multicollinearity really occurs. Heatmaps are great to detect this kind of situations and in problems dominated by feature selection, like ours, they are an essential tool.

Another thing that got my attention was the 'SalePrice' correlations. We can see our well-known 'GrLivArea', 'TotalBsmtSF', and 'OverallQual' saying a big 'Hi!', but we can also see many other variables that should be taken into account. That's what we will do next.

#### 'SalePrice' correlation matrix (zoomed heatmap style)

```
In [12]: #saleprice correlation matrix
k = 10 #number of variables for heatmap
cols = corrmat.nlargest(k, 'SalePrice')['SalePrice'].index
cm = np.corrcoef(df_train[cols].values.T)
sns.set(font_scale=1.25)
hm = sns.heatmap(cm, cbar=True, annot=True, square=True, fmt='.2f', annot_kw
s={'size': 10}, yticklabels=cols.values, xticklabels=cols.values)
plt.show()
```



According to our crystal ball, these are the variables most correlated with 'SalePrice'. My thoughts on this:

- 'OverallQual', 'GrLivArea' and 'TotalBsmtSF' are strongly correlated with 'SalePrice'. Check!
- 'GarageCars' and 'GarageArea' are also some of the most strongly correlated variables. However, as we discussed in the last sub-point, the number of cars that fit into the garage is a consequence of the garage area. 'GarageCars' and 'GarageArea' are like twin brothers. You'll never be able to distinguish them. Therefore, we just need one of these variables in our analysis (we can keep 'GarageCars' since its correlation with 'SalePrice' is higher).
- 'TotalBsmtSF' and '1stFloor' also seem to be twin brothers. We can keep 'TotalBsmtSF' just to say that our first guess was right (re-read 'So... What can we expect?').
- 'FullBath'?? Really?
- 'TotRmsAbvGrd' and 'GrLivArea', twin brothers again. Is this dataset from Chernobyl?
- Ah... 'YearBuilt'... It seems that 'YearBuilt' is slightly correlated with 'SalePrice'. Honestly, it scares me to think about 'YearBuilt' because I start feeling that we should do a little bit of time-series analysis to get this right. I'll leave this as a homework for you.

Let's proceed to the scatter plots.

### Scatter plots between 'SalePrice' and correlated variables (move like Jagger style)

Get ready for what you're about to see. I must confess that the first time I saw these scatter plots I was totally blown away! So much information in so short space... It's just amazing. Once more, thank you @seaborn! You make me 'move like Jagger'!

```
In [13]: #scatterplot
sns.set()
cols = ['SalePrice', 'OverallQual', 'GrLivArea', 'GarageCars', 'TotalBsmtSF',
        'FullBath', 'YearBuilt']
sns.pairplot(df_train[cols], size = 2.5)
plt.show();
```



Although we already know some of the main figures, this mega scatter plot gives us a reasonable idea about variables relationships.

One of the figures we may find interesting is the one between 'TotalBsmtSF' and 'GrLiveArea'. In this figure we can see the dots drawing a linear line, which almost acts like a border. It totally makes sense that the majority of the dots stay below that line. Basement areas can be equal to the above ground living area, but it is not expected a basement area bigger than the above ground living area (unless you're trying to buy a bunker).

The plot concerning 'SalePrice' and 'YearBuilt' can also make us think. In the bottom of the 'dots cloud', we see what almost appears to be a shy exponential function (be creative). We can also see this same tendency in the upper limit of the 'dots cloud' (be even more creative). Also, notice how the set of dots regarding the last years tend to stay above this limit (I just wanted to say that prices are increasing faster now).

Ok, enough of Rorschach test for now. Let's move forward to what's missing: missing data!

## 4. Missing data

Important questions when thinking about missing data:

- How prevalent is the missing data?
- Is missing data random or does it have a pattern?

The answer to these questions is important for practical reasons because missing data can imply a reduction of the sample size. This can prevent us from proceeding with the analysis. Moreover, from a substantive perspective, we need to ensure that the missing data process is not biased and hiding an inconvenient truth.

```
In [14]: #missing data
total = df_train.isnull().sum().sort_values(ascending=False)
percent = (df_train.isnull().sum()/df_train.isnull().count()).sort_values(ascending=False)
missing_data = pd.concat([total, percent], axis=1, keys=['Total', 'Percent'])
missing_data.head(20)
```

Out[14]:

	Total	Percent
<b>PoolQC</b>	1453	0.995205
<b>MiscFeature</b>	1406	0.963014
<b>Alley</b>	1369	0.937671
<b>Fence</b>	1179	0.807534
<b>FireplaceQu</b>	690	0.472603
<b>LotFrontage</b>	259	0.177397
<b>GarageCond</b>	81	0.055479
<b>GarageType</b>	81	0.055479
<b>GarageYrBlt</b>	81	0.055479
<b>GarageFinish</b>	81	0.055479
<b>GarageQual</b>	81	0.055479
<b>BsmtExposure</b>	38	0.026027
<b>BsmtFinType2</b>	38	0.026027
<b>BsmtFinType1</b>	37	0.025342
<b>BsmtCond</b>	37	0.025342
<b>BsmtQual</b>	37	0.025342
<b>MasVnrArea</b>	8	0.005479
<b>MasVnrType</b>	8	0.005479
<b>Electrical</b>	1	0.000685
<b>Utilities</b>	0	0.000000

Let's analyse this to understand how to handle the missing data.

We'll consider that when more than 15% of the data is missing, we should delete the corresponding variable and pretend it never existed. This means that we will not try any trick to fill the missing data in these cases. According to this, there is a set of variables (e.g. 'PoolQC', 'MiscFeature', 'Alley', etc.) that we should delete. The point is: will we miss this data? I don't think so. None of these variables seem to be very important, since most of them are not aspects in which we think about when buying a house (maybe that's the reason why data is missing?). Moreover, looking closer at the variables, we could say that variables like 'PoolQC', 'MiscFeature' and 'FireplaceQu' are strong candidates for outliers, so we'll be happy to delete them.

In what concerns the remaining cases, we can see that 'GarageX' variables have the same number of missing data. I bet missing data refers to the same set of observations (although I will not check it; it's just 5% and we should not spend 20 minutes on problems). Since the most important information regarding garages is expressed by 'GarageCars' and considering that we are just talking about 5% of missing data, I'll delete the mentioned 'GarageX' variables. The same logic applies to 'BsmtX' variables.

Regarding 'MasVnrArea' and 'MasVnrType', we can consider that these variables are not essential. Furthermore, they have a strong correlation with 'YearBuilt' and 'OverallQual' which are already considered. Thus, we will not lose information if we delete 'MasVnrArea' and 'MasVnrType'.

Finally, we have one missing observation in 'Electrical'. Since it is just one observation, we'll delete this observation and keep the variable.

In summary, to handle missing data, we'll delete all the variables with missing data, except the variable 'Electrical'. In 'Electrical' we'll just delete the observation with missing data.

```
In [15]: #dealing with missing data
df_train = df_train.drop((missing_data[missing_data['Total'] > 1]).index,1)
df_train = df_train.drop(df_train.loc[df_train['Electrical'].isnull()].index)
df_train.isnull().sum().max() #just checking that there's no missing data missing...
```

Out[15]: 0

## Out liars!

Outliers is also something that we should be aware of. Why? Because outliers can markedly affect our models and can be a valuable source of information, providing us insights about specific behaviours.

Outliers is a complex subject and it deserves more attention. Here, we'll just do a quick analysis through the standard deviation of 'SalePrice' and a set of scatter plots.

## Univariate analysis

The primary concern here is to establish a threshold that defines an observation as an outlier. To do so, we'll standardize the data. In this context, data standardization means converting data values to have mean of 0 and a standard deviation of 1.

```
In [16]: #standardizing data
saleprice_scaled = StandardScaler().fit_transform(df_train['SalePrice'][:,n
p.newaxis]);
low_range = saleprice_scaled[saleprice_scaled[:,0].argsort()][:10]
high_range= saleprice_scaled[saleprice_scaled[:,0].argsort()][-10:]
print('outer range (low) of the distribution:')
print(low_range)
print('\nouter range (high) of the distribution:')
print(high_range)

outer range (low) of the distribution:
[[-1.83820775]
 [-1.83303414]
 [-1.80044422]
 [-1.78282123]
 [-1.77400974]
 [-1.62295562]
 [-1.6166617 ]
 [-1.58519209]
 [-1.58519209]
 [-1.57269236]]

outer range (high) of the distribution:
[[3.82758058]
 [4.0395221 ]
 [4.49473628]
 [4.70872962]
 [4.728631  ]
 [5.06034585]
 [5.42191907]
 [5.58987866]
 [7.10041987]
 [7.22629831]]
```

How 'SalePrice' looks with her new clothes:

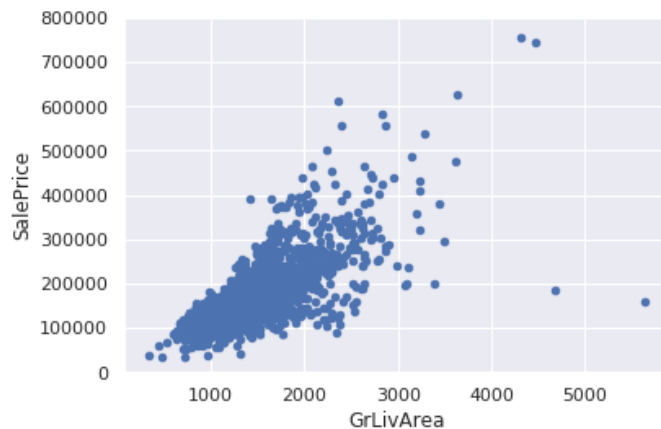
- Low range values are similar and not too far from 0.
- High range values are far from 0 and the 7.something values are really out of range.

For now, we'll not consider any of these values as an outlier but we should be careful with those two 7.something values.

## Bivariate analysis

We already know the following scatter plots by heart. However, when we look to things from a new perspective, there's always something to discover. As Alan Kay said, 'a change in perspective is worth 80 IQ points'.

```
In [17]: #bivariate analysis saleprice/grlivarea
var = 'GrLivArea'
data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
data.plot.scatter(x=var, y='SalePrice', ylim=(0,800000));
```

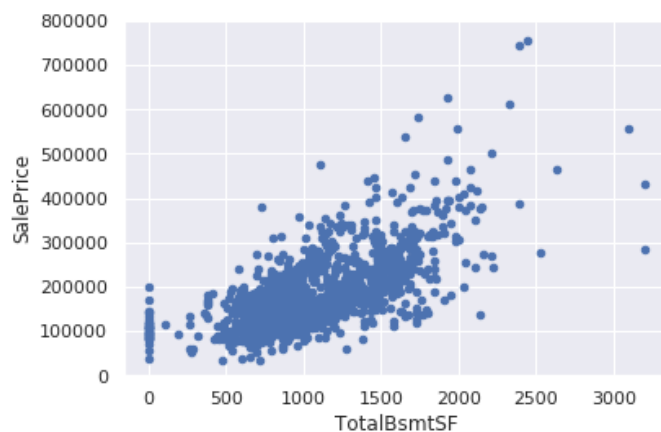


What has been revealed:

- The two values with bigger 'GrLivArea' seem strange and they are not following the crowd. We can speculate why this is happening. Maybe they refer to agricultural area and that could explain the low price. I'm not sure about this but I'm quite confident that these two points are not representative of the typical case. Therefore, we'll define them as outliers and delete them.
- The two observations in the top of the plot are those 7.something observations that we said we should be careful about. They look like two special cases, however they seem to be following the trend. For that reason, we will keep them.

```
In [18]: #deleting points
df_train.sort_values(by = 'GrLivArea', ascending = False)[:2]
df_train = df_train.drop(df_train[df_train['Id'] == 1299].index)
df_train = df_train.drop(df_train[df_train['Id'] == 524].index)
```

```
In [19]: #bivariate analysis saleprice/grlivarea
var = 'TotalBsmtSF'
data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
data.plot.scatter(x=var, y='SalePrice', ylim=(0,800000));
```



We can feel tempted to eliminate some observations (e.g. TotalBsmtSF > 3000) but I suppose it's not worth it. We can live with that, so we'll not do anything.

## 5. Getting hard core

In Ayn Rand's novel, 'Atlas Shrugged', there is an often-repeated question: who is John Galt? A big part of the book is about the quest to discover the answer to this question.

I feel Randian now. Who is 'SalePrice'?

The answer to this question lies in testing for the assumptions underlying the statistical bases for multivariate analysis. We already did some data cleaning and discovered a lot about 'SalePrice'. Now it's time to go deep and understand how 'SalePrice' complies with the statistical assumptions that enables us to apply multivariate techniques.

According to [Hair et al. \(2013\)](https://amzn.to/2uC3j9p) (<https://amzn.to/2uC3j9p>), four assumptions should be tested:

- **Normality** - When we talk about normality what we mean is that the data should look like a normal distribution. This is important because several statistic tests rely on this (e.g. t-statistics). In this exercise we'll just check univariate normality for 'SalePrice' (which is a limited approach). Remember that univariate normality doesn't ensure multivariate normality (which is what we would like to have), but it helps. Another detail to take into account is that in big samples (>200 observations) normality is not such an issue. However, if we solve normality, we avoid a lot of other problems (e.g. heteroscedacity) so that's the main reason why we are doing this analysis.
- **Homoscedasticity** - I just hope I wrote it right. Homoscedasticity refers to the 'assumption that dependent variable(s) exhibit equal levels of variance across the range of predictor variable(s)' ([Hair et al., 2013](https://amzn.to/2uC3j9p)) (<https://amzn.to/2uC3j9p>). Homoscedasticity is desirable because we want the error term to be the same across all values of the independent variables.
- **Linearity** - The most common way to assess linearity is to examine scatter plots and search for linear patterns. If patterns are not linear, it would be worthwhile to explore data transformations. However, we'll not get into this because most of the scatter plots we've seen appear to have linear relationships.
- **Absence of correlated errors** - Correlated errors, like the definition suggests, happen when one error is correlated to another. For instance, if one positive error makes a negative error systematically, it means that there's a relationship between these variables. This occurs often in time series, where some patterns are time related. We'll also not get into this. However, if you detect something, try to add a variable that can explain the effect you're getting. That's the most common solution for correlated errors.

What do you think Elvis would say about this long explanation? 'A little less conversation, a little more action please'? Probably... By the way, do you know what was Elvis's last great hit?

(...)

The bathroom floor.

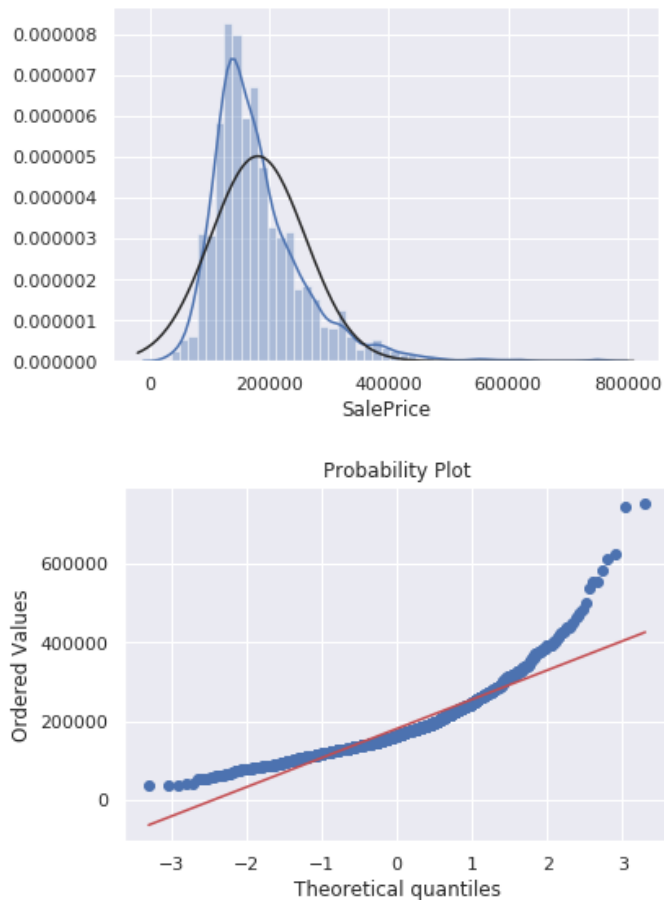
### In the search for normality

The point here is to test 'SalePrice' in a very lean way. We'll do this paying attention to:

- **Histogram** - Kurtosis and skewness.
- **Normal probability plot** - Data distribution should closely follow the diagonal that represents the normal distribution.



```
In [20]: #histogram and normal probability plot
sns.distplot(df_train['SalePrice'], fit=norm);
fig = plt.figure()
res = stats.probplot(df_train['SalePrice'], plot=plt)
```



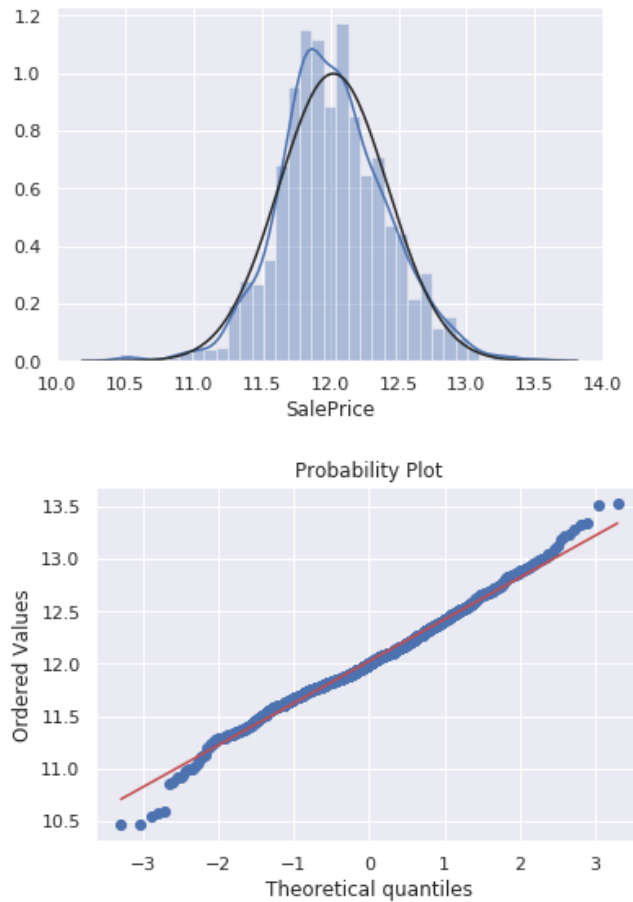
Ok, 'SalePrice' is not normal. It shows 'peakedness', positive skewness and does not follow the diagonal line.

But everything's not lost. A simple data transformation can solve the problem. This is one of the awesome things you can learn in statistical books: in case of positive skewness, log transformations usually works well. When I discovered this, I felt like an Hogwarts' student discovering a new cool spell.

*Avada kedavra!*

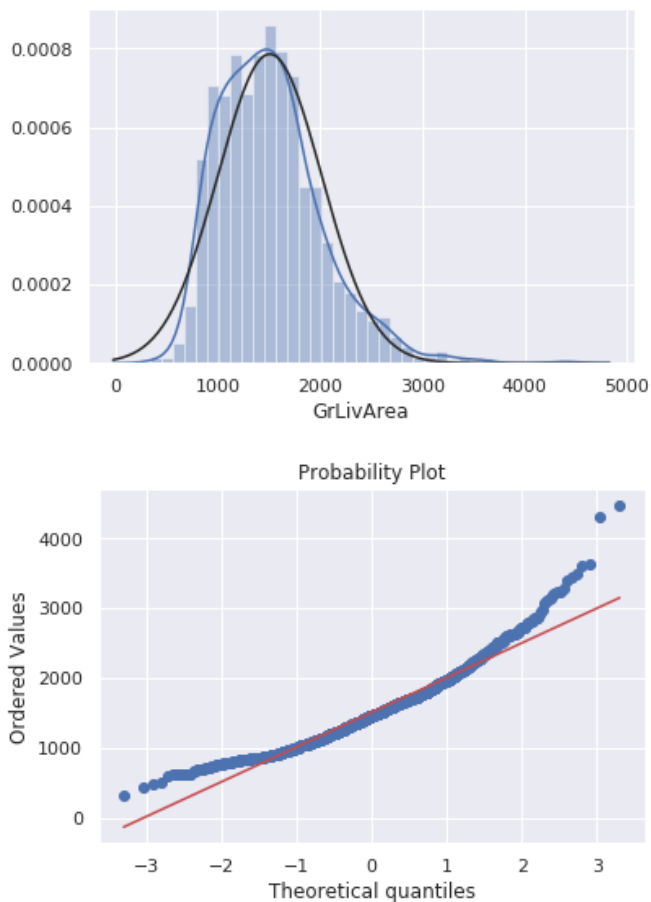
```
In [21]: #applying log transformation
df_train['SalePrice'] = np.log(df_train['SalePrice'])
```

```
In [22]: #transformed histogram and normal probability plot  
sns.distplot(df_train['SalePrice'], fit=norm);  
fig = plt.figure()  
res = stats.probplot(df_train['SalePrice'], plot=plt)
```



Done! Let's check what's going on with 'GrLivArea'.

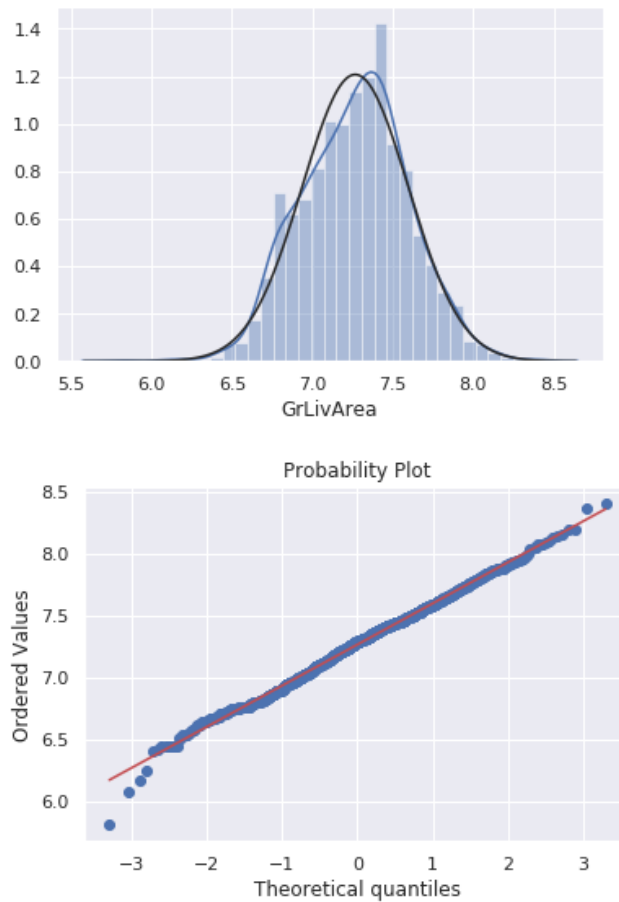
```
In [23]: #histogram and normal probability plot
sns.distplot(df_train['GrLivArea'], fit=norm);
fig = plt.figure()
res = stats.probplot(df_train['GrLivArea'], plot=plt)
```



Tastes like skewness... *Avada kedavra!*

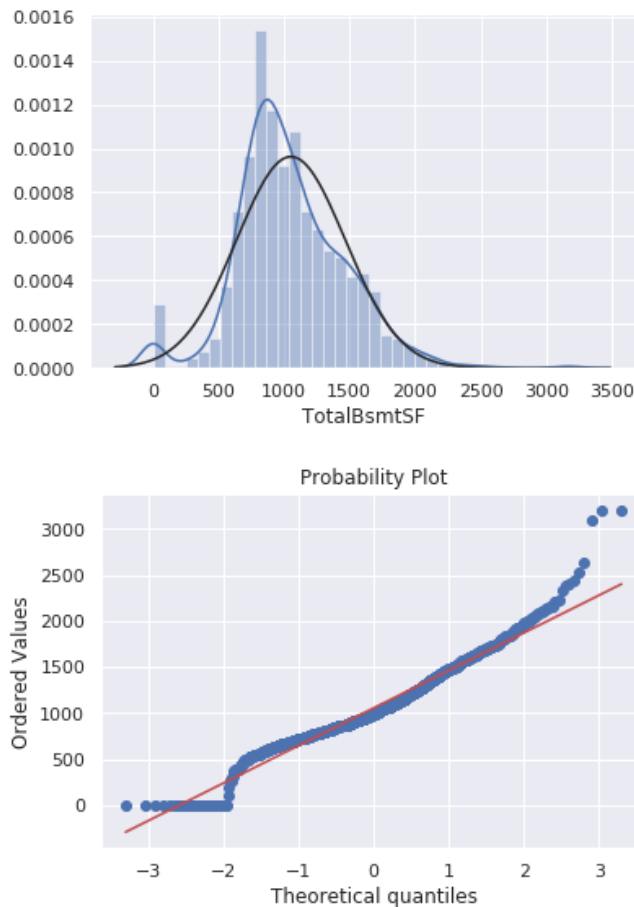
```
In [24]: #data transformation
df_train['GrLivArea'] = np.log(df_train['GrLivArea'])
```

```
In [25]: #transformed histogram and normal probability plot
sns.distplot(df_train['GrLivArea'], fit=norm);
fig = plt.figure()
res = stats.probplot(df_train['GrLivArea'], plot=plt)
```



Next, please...

```
In [26]: #histogram and normal probability plot
sns.distplot(df_train['TotalBsmtSF'], fit=norm);
fig = plt.figure()
res = stats.probplot(df_train['TotalBsmtSF'], plot=plt)
```



Ok, now we are dealing with the big boss. What do we have here?

- Something that, in general, presents skewness.
- A significant number of observations with value zero (houses without basement).
- A big problem because the value zero doesn't allow us to do log transformations.

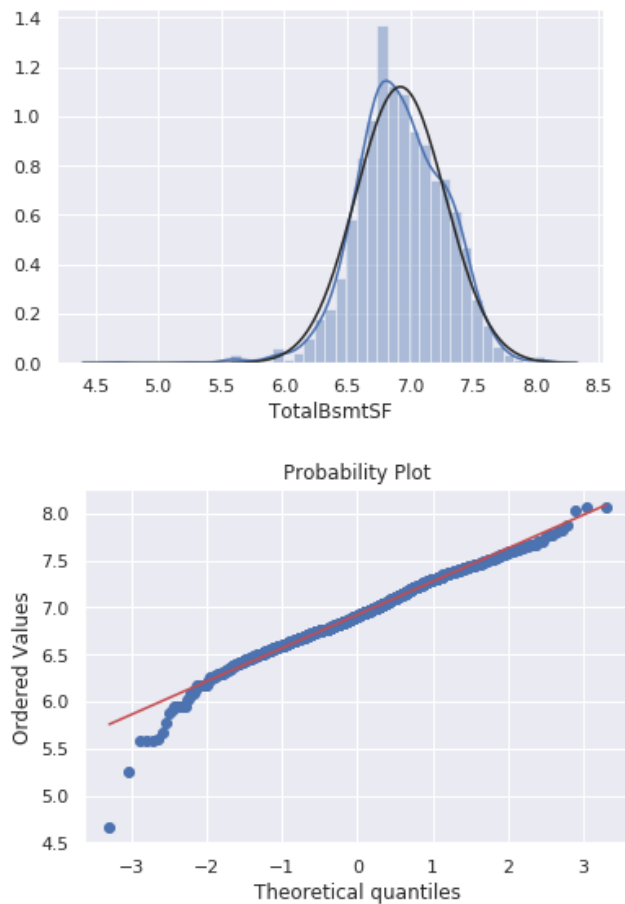
To apply a log transformation here, we'll create a variable that can get the effect of having or not having basement (binary variable). Then, we'll do a log transformation to all the non-zero observations, ignoring those with value zero. This way we can transform data, without losing the effect of having or not basement.

I'm not sure if this approach is correct. It just seemed right to me. That's what I call 'high risk engineering'.

```
In [27]: #create column for new variable (one is enough because it's a binary categorical feature)
#if area>0 it gets 1, for area==0 it gets 0
df_train['HasBsmt'] = pd.Series(len(df_train['TotalBsmtSF']), index=df_train.index)
df_train['HasBsmt'] = 0
df_train.loc[df_train['TotalBsmtSF']>0, 'HasBsmt'] = 1
```

```
In [28]: #transform data
df_train.loc[df_train['HasBsmt']==1, 'TotalBsmtSF'] = np.log(df_train['TotalBsmtSF'])
```

```
In [29]: #histogram and normal probability plot
sns.distplot(df_train[df_train['TotalBsmtSF']>0]['TotalBsmtSF'], fit=norm);
fig = plt.figure()
res = stats.probplot(df_train[df_train['TotalBsmtSF']>0]['TotalBsmtSF'], plo
t=plt)
```

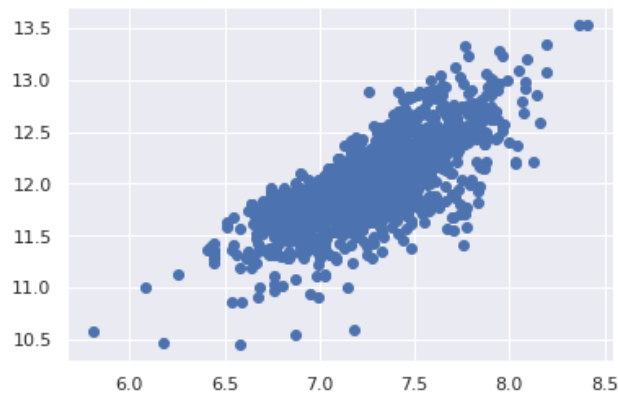


### In the search for writing 'homoscedasticity' right at the first attempt

The best approach to test homoscedasticity for two metric variables is graphically. Departures from an equal dispersion are shown by such shapes as cones (small dispersion at one side of the graph, large dispersion at the opposite side) or diamonds (a large number of points at the center of the distribution).

Starting by 'SalePrice' and 'GrLivArea'...

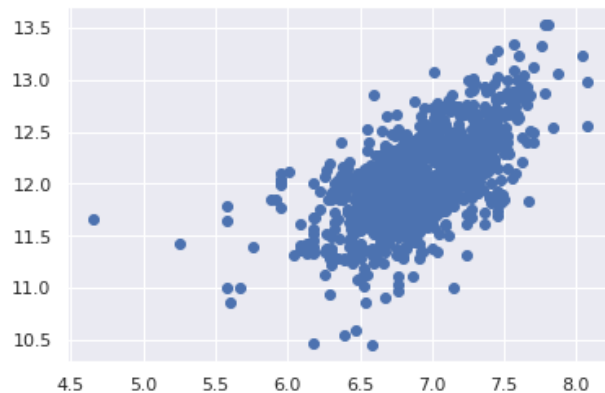
```
In [30]: #scatter plot
plt.scatter(df_train['GrLivArea'], df_train['SalePrice']);
```



Older versions of this scatter plot (previous to log transformations), had a conic shape (go back and check 'Scatter plots between 'SalePrice' and correlated variables (move like Jagger style)'). As you can see, the current scatter plot doesn't have a conic shape anymore. That's the power of normality! Just by ensuring normality in some variables, we solved the homoscedasticity problem.

Now let's check 'SalePrice' with 'TotalBsmtSF'.

```
In [31]: #scatter plot
plt.scatter(df_train[df_train['TotalBsmtSF']>0]['TotalBsmtSF'], df_train[df_train['TotalBsmtSF']>0]['SalePrice']);
```



We can say that, in general, 'SalePrice' exhibit equal levels of variance across the range of 'TotalBsmtSF'. Cool!

## Last but not the least, dummy variables

Easy mode.

```
In [32]: #convert categorical variable into dummy
df_train = pd.get_dummies(df_train)
```

## Conclusion

That's it! We reached the end of our exercise.

Throughout this kernel we put in practice many of the strategies proposed by [Hair et al. \(2013\)](https://amzn.to/2uC3j9p). We philosophied about the variables, we analysed 'SalePrice' alone and with the most correlated variables, we dealt with missing data and outliers, we tested some of the fundamental statistical assumptions and we even transformed categorial variables into dummy variables. That's a lot of work that Python helped us make easier.

But the quest is not over. Remember that our story stopped in the Facebook research. Now it's time to give a call to 'SalePrice' and invite her to dinner. Try to predict her behaviour. Do you think she's a girl that enjoys regularized linear regression approaches? Or do you think she prefers ensemble methods? Or maybe something else?

It's up to you to find out.

## References

- [My blog](http://pmarcelino.com) (<http://pmarcelino.com>)
- [My other kernels](https://www.kaggle.com/pmarcelino/data-analysis-and-feature-extraction-with-python) (<https://www.kaggle.com/pmarcelino/data-analysis-and-feature-extraction-with-python>)
- [Hair et al., 2013, Multivariate Data Analysis, 7th Edition](https://amzn.to/2JuDmvo) (<https://amzn.to/2JuDmvo>)

## Acknowledgements

Thanks to [João Rico](https://www.linkedin.com/in/joaomiguelrico/) (<https://www.linkedin.com/in/joaomiguelrico/>) for reading drafts of this.



## Trying out a linear model:

Author: Alexandru Papiu (@apapiu (<https://twitter.com/apapiu>), [GitHub \(https://github.com/apapiu\)](https://github.com/apapiu))

If you use parts of this notebook in your own scripts, please give some sort of credit (for example link back to this).  
Thanks!

There have been a few [great \(https://www.kaggle.com/comartel/house-prices-advanced-regression-techniques/house-price-xgboost-starter/run/348739\)](https://www.kaggle.com/comartel/house-prices-advanced-regression-techniques/house-price-xgboost-starter/run/348739) [scripts \(https://www.kaggle.com/zoupet/house-prices-advanced-regression-techniques/xgboost-10-kfolds-with-scikit-learn/run/357561\)](https://www.kaggle.com/zoupet/house-prices-advanced-regression-techniques/xgboost-10-kfolds-with-scikit-learn/run/357561) on [xgboost \(https://www.kaggle.com/tadepalli/house-prices-advanced-regression-techniques/xgboost-with-n-trees-autostop-0-12638/run/353049\)](https://www.kaggle.com/tadepalli/house-prices-advanced-regression-techniques/xgboost-with-n-trees-autostop-0-12638/run/353049) already so I'd figured I'd try something simpler: a regularized linear regression model. Surprisingly it does really well with very little feature engineering. The key point is to log\_transform the numeric variables since most of them are skewed.

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib

import matplotlib.pyplot as plt
from scipy.stats import skew
from scipy.stats.stats import pearsonr

%config InlineBackend.figure_format = 'retina' #set 'png' here when working
on notebook
%matplotlib inline
```

```
In [2]: train = pd.read_csv("../input/train.csv")
test = pd.read_csv("../input/test.csv")
```

```
In [3]: train.head()
```

```
Out[3]:
```

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

5 rows × 81 columns

```
In [4]: all_data = pd.concat((train.loc[:, 'MSSubClass': 'SaleCondition'],
                             test.loc[:, 'MSSubClass': 'SaleCondition']))
```

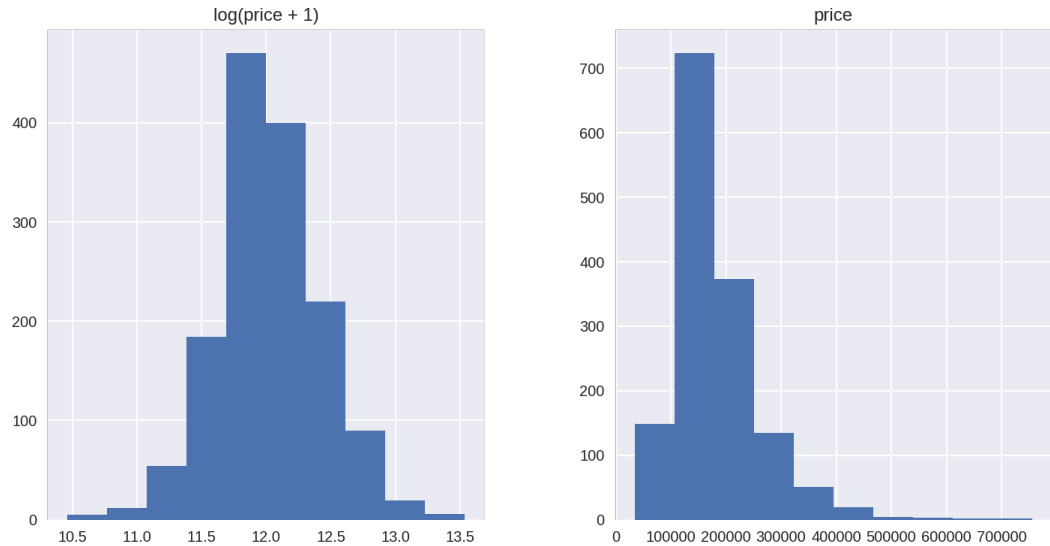
## Data preprocessing:

We're not going to do anything fancy here:

- First I'll transform the skewed numeric features by taking  $\log(\text{feature} + 1)$  - this will make the features more normal
- Create Dummy variables for the categorical features
- Replace the numeric missing values (NaN's) with the mean of their respective columns

```
In [5]: matplotlib.rcParams['figure.figsize'] = (12.0, 6.0)
prices = pd.DataFrame({"price":train["SalePrice"], "log(price + 1)":np.log1p(
train["SalePrice"])})
prices.hist()
```

```
Out[5]: array([[<matplotlib.axes._subplots.AxesSubplot object at 0x7f636aca7978>,
<matplotlib.axes._subplots.AxesSubplot object at 0x7f63678dc240>]], dtype=object)
```



```
In [6]: #log transform the target:
train["SalePrice"] = np.log1p(train["SalePrice"])

#log transform skewed numeric features:
numeric_feats = all_data.dtypes[all_data.dtypes != "object"].index

skewed_feats = train[numeric_feats].apply(lambda x: skew(x.dropna())) #compute skewness
skewed_feats = skewed_feats[skewed_feats > 0.75]
skewed_feats = skewed_feats.index

all_data[skewed_feats] = np.log1p(all_data[skewed_feats])
```

```
In [7]: all_data = pd.get_dummies(all_data)
```

```
In [8]: #filling NA's with the mean of the column:
all_data = all_data.fillna(all_data.mean())
```

```
In [9]: #creating matrices for sklearn:
X_train = all_data[:train.shape[0]]
X_test = all_data[train.shape[0]:]
y = train.SalePrice
```

## Models

Now we are going to use regularized linear regression models from the scikit learn module. I'm going to try both `l_1`(Lasso) and `l_2`(Ridge) regularization. I'll also define a function that returns the cross-validation rmse error so we can evaluate our models and pick the best tuning par

```
In [10]: from sklearn.linear_model import Ridge, RidgeCV, ElasticNet, LassoCV, LassoLarsCV
from sklearn.model_selection import cross_val_score

def rmse_cv(model):
    rmse= np.sqrt(-cross_val_score(model, X_train, y, scoring="neg_mean_squared_error", cv = 5))
    return(rmse)
```

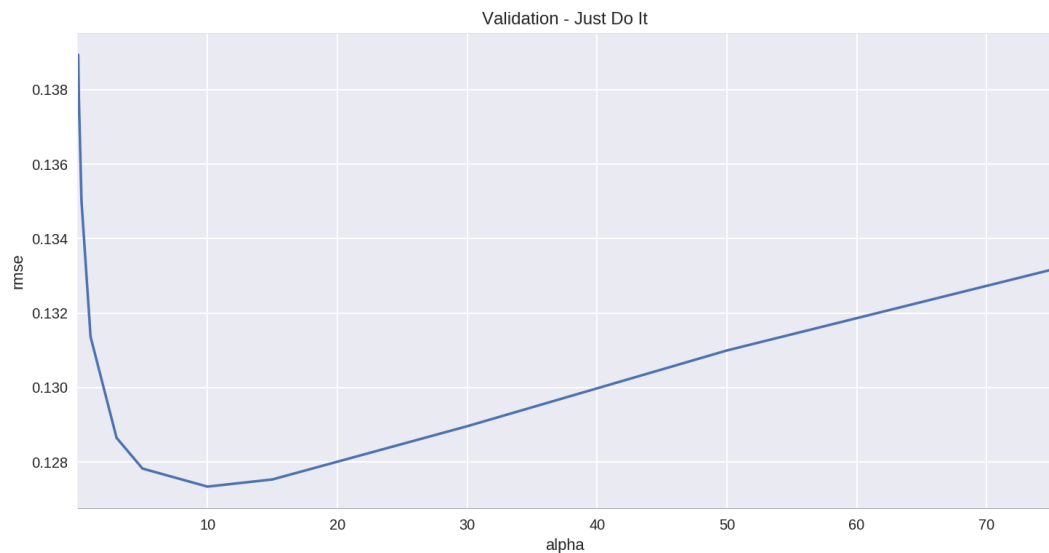
```
In [11]: model_ridge = Ridge()
```

The main tuning parameter for the Ridge model is alpha - a regularization parameter that measures how flexible our model is. The higher the regularization the less prone our model will be to overfit. However it will also lose flexibility and might not capture all of the signal in the data.

```
In [12]: alphas = [0.05, 0.1, 0.3, 1, 3, 5, 10, 15, 30, 50, 75]
cv_ridge = [rmse_cv(Ridge(alpha = alpha)).mean()
             for alpha in alphas]
```

```
In [13]: cv_ridge = pd.Series(cv_ridge, index = alphas)
cv_ridge.plot(title = "Validation - Just Do It")
plt.xlabel("alpha")
plt.ylabel("rmse")
```

```
Out[13]: <matplotlib.text.Text at 0x7f6364300eb8>
```



Note the U-ish shaped curve above. When alpha is too large the regularization is too strong and the model cannot capture all the complexities in the data. If however we let the model be too flexible (alpha small) the model begins to overfit. A value of alpha = 10 is about right based on the plot above.

```
In [14]: cv_ridge.min()
```

```
Out[14]: 0.12733734668670765
```

So for the Ridge regression we get a rmse of about 0.127

Let's try out the Lasso model. We will do a slightly different approach here and use the built in Lasso CV to figure out the best alpha for us. For some reason the alphas in Lasso CV are really the inverse or the alphas in Ridge.

```
In [15]: model_lasso = LassoCV(alphas = [1, 0.1, 0.001, 0.0005]).fit(X_train, y)
```

```
In [16]: rmse_cv(model_lasso).mean()
```

```
Out[16]: 0.12314421090977427
```

Nice! The lasso performs even better so we'll just use this one to predict on the test set. Another neat thing about the Lasso is that it does feature selection for you - setting coefficients of features it deems unimportant to zero. Let's take a look at the coefficients:

```
In [17]: coef = pd.Series(model_lasso.coef_, index = X_train.columns)
```

```
In [18]: print("Lasso picked " + str(sum(coef != 0)) + " variables and eliminated the  
other " + str(sum(coef == 0)) + " variables")
```

```
Lasso picked 110 variables and eliminated the other 178 variables
```

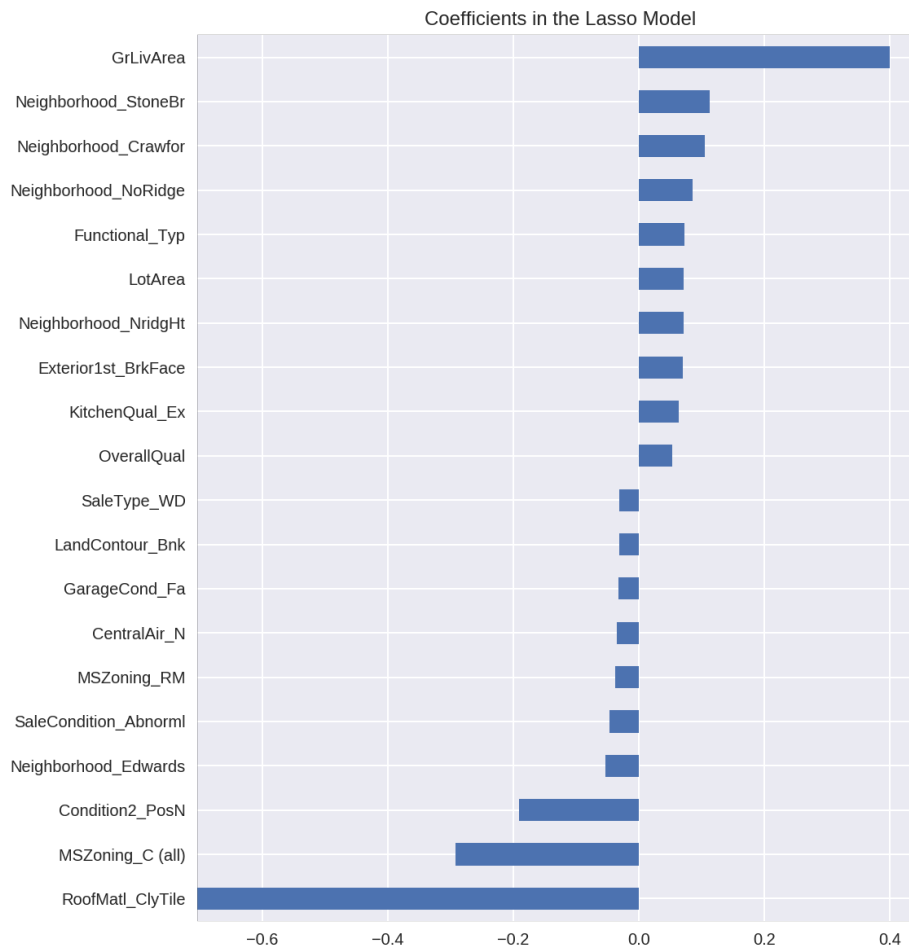
Good job Lasso. One thing to note here however is that the features selected are not necessarily the "correct" ones - especially since there are a lot of collinear features in this dataset. One idea to try here is run Lasso a few times on bootstrapped samples and see how stable the feature selection is.

We can also take a look directly at what the most important coefficients are:

```
In [19]: imp_coef = pd.concat([coef.sort_values().head(10),  
                             coef.sort_values().tail(10)])
```

```
In [20]: matplotlib.rcParams['figure.figsize'] = (8.0, 10.0)
imp_coef.plot(kind = "barh")
plt.title("Coefficients in the Lasso Model")
```

```
Out[20]: <matplotlib.text.Text at 0x7f63605ca4e0>
```



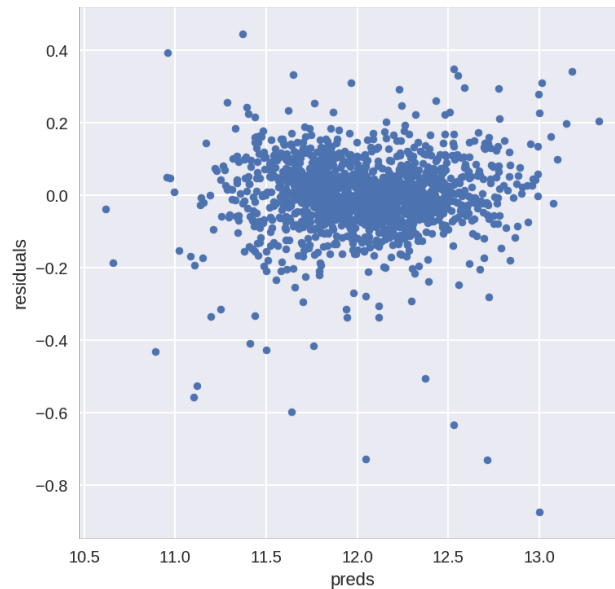
The most important positive feature is `GrLivArea` - the above ground area by area square feet. This definitely sense. Then a few other location and quality features contributed positively. Some of the negative features make less sense and would be worth looking into more - it seems like they might come from unbalanced categorical variables.

Also note that unlike the feature importance you'd get from a random forest these are *actual* coefficients in your model - so you can say precisely why the predicted price is what it is. The only issue here is that we `log_transformed` both the target and the numeric features so the actual magnitudes are a bit hard to interpret.

```
In [21]: #let's look at the residuals as well:
matplotlib.rcParams['figure.figsize'] = (6.0, 6.0)

preds = pd.DataFrame({"preds":model_lasso.predict(X_train), "true":y})
preds["residuals"] = preds["true"] - preds["preds"]
preds.plot(x = "preds", y = "residuals",kind = "scatter")
```

Out[21]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7f63604cb160>



The residual plot looks pretty good. To wrap it up let's predict on the test set and submit on the leaderboard:

### Adding an xgboost model:

Let's add an xgboost model to our linear model to see if we can improve our score:

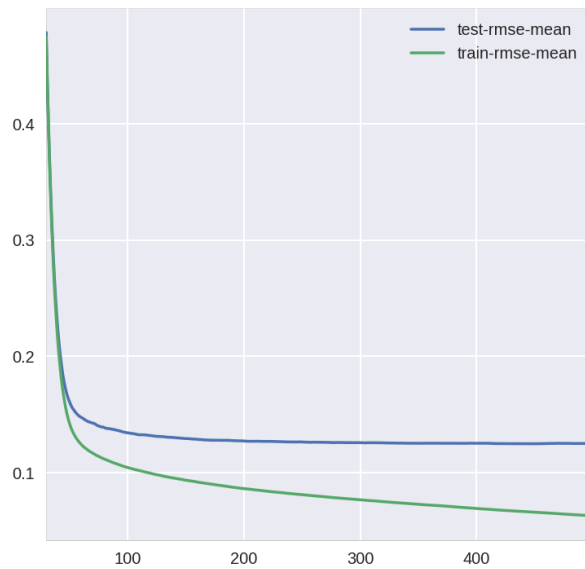
```
In [22]: import xgboost as xgb
```

```
In [23]: dtrain = xgb.DMatrix(X_train, label = y)
dtest = xgb.DMatrix(X_test)

params = {"max_depth":2, "eta":0.1}
model = xgb.cv(params, dtrain, num_boost_round=500, early_stopping_rounds=100)
```

```
In [24]: model.loc[30:,[ "test-rmse-mean", "train-rmse-mean"]].plot()
```

```
Out[24]: <matplotlib.axes._subplots.AxesSubplot at 0x7f6360465d68>
```



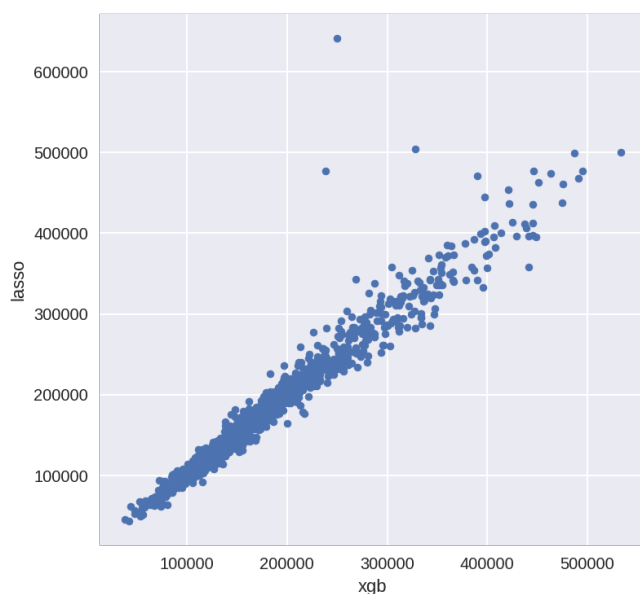
```
In [25]: model_xgb = xgb.XGBRegressor(n_estimators=360, max_depth=2, learning_rate=0.1)
#the params were tuned using xgb.cv
model_xgb.fit(X_train, y)
```

```
Out[25]: XGBRegressor(base_score=0.5, colsample_bylevel=1, colsample_bytree=1, gamma=0,
learning_rate=0.1, max_delta_step=0, max_depth=2,
min_child_weight=1, missing=None, n_estimators=360, nthread=-1,
objective='reg:linear', reg_alpha=0, reg_lambda=1,
scale_pos_weight=1, seed=0, silent=True, subsample=1)
```

```
In [26]: xgb_preds = np.expml(model_xgb.predict(X_test))
lasso_preds = np.expml(model_lasso.predict(X_test))
```

```
In [27]: predictions = pd.DataFrame({"xgb":xgb_preds, "lasso":lasso_preds})
         predictions.plot(x = "xgb", y = "lasso", kind = "scatter")
```

```
Out[27]: <matplotlib.axes._subplots.AxesSubplot at 0x7f6360469f60>
```



Many times it makes sense to take a weighted average of uncorrelated results - this usually improves the score although in this case it doesn't help that much.

```
In [28]: preds = 0.7*lasso_preds + 0.3*xgb_preds
```

```
In [29]: solution = pd.DataFrame({"id":test.Id, "SalePrice":preds})
         solution.to_csv("ridge_sol.csv", index = False)
```

## Trying out keras?

Feedforward Neural Nets doesn't seem to work well at all...I wonder why.

```
In [30]: from keras.layers import Dense
         from keras.models import Sequential
         from keras.regularizers import l1
         from sklearn.preprocessing import StandardScaler
         from sklearn.model_selection import train_test_split
```

Using TensorFlow backend.

```
In [31]: X_train = StandardScaler().fit_transform(X_train)
```

```
In [32]: X_tr, X_val, y_tr, y_val = train_test_split(X_train, y, random_state = 3)
```

```
In [33]: X_tr.shape
```

```
Out[33]: (1095, 288)
```



In [34]: X\_tr

```
Out[34]: array([[ 1.00573733,  0.68066137, -0.46001991, ..., -0.11785113,
                  0.4676514 , -0.30599503],
                [-1.12520184,  0.60296111,  0.03113183, ..., -0.11785113,
                  0.4676514 , -0.30599503],
                [-1.12520184, -0.02865265, -0.74027492, ..., -0.11785113,
                  0.4676514 , -0.30599503],
                ...,
                [ 0.16426234, -0.87075036, -0.81954431, ..., -0.11785113,
                  -2.13834494, -0.30599503],
                [ 0.92361154, -0.30038284, -0.44275864, ..., -0.11785113,
                  0.4676514 , -0.30599503],
                [ 0.83656519,  1.98505948,  0.46455838, ..., -0.11785113,
                  0.4676514 , -0.30599503]])
```

```
In [35]: model = Sequential()
          #model.add(Dense(256, activation="relu", input_dim = X_train.shape[1]))
          model.add(Dense(1, input_dim = X_train.shape[1], W_regularizer=l1(0.001)))

          model.compile(loss = "mse", optimizer = "adam")
```

```
/opt/conda/lib/python3.6/site-packages/ipykernel/__main__.py:3: UserWarning:
Update your `Dense` call to the Keras 2 API: `Dense(1, input_dim=288, kernel_
regularizer=<keras.reg...)`
  app.launch_new_instance()
```

In [36]: model.summary()

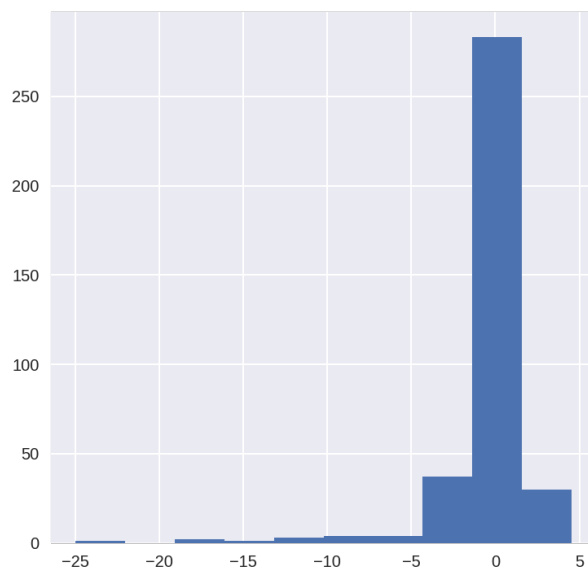
Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 1)	289
Total params: 289		
Trainable params: 289		
Non-trainable params: 0		

```
In [37]: hist = model.fit(X_tr, y_tr, validation_data = (X_val, y_val))
```

```
Train on 1095 samples, validate on 365 samples
Epoch 1/10
1095/1095 [=====] - 0s - loss: 147.0313 - val_loss: 149.9484
Epoch 2/10
1095/1095 [=====] - 0s - loss: 144.6278 - val_loss: 150.5536
Epoch 3/10
1095/1095 [=====] - 0s - loss: 142.8769 - val_loss: 151.4110
Epoch 4/10
1095/1095 [=====] - 0s - loss: 141.3150 - val_loss: 152.4455
Epoch 5/10
1095/1095 [=====] - 0s - loss: 139.8918 - val_loss: 153.7267
Epoch 6/10
1095/1095 [=====] - 0s - loss: 138.4870 - val_loss: 155.0268
Epoch 7/10
1095/1095 [=====] - 0s - loss: 137.1905 - val_loss: 156.3740
Epoch 8/10
1095/1095 [=====] - 0s - loss: 135.9163 - val_loss: 157.8431
Epoch 9/10
1095/1095 [=====] - 0s - loss: 134.6428 - val_loss: 159.4450
Epoch 10/10
1095/1095 [=====] - 0s - loss: 133.3451 - val_loss: 161.2576
```

```
In [38]: pd.Series(model.predict(X_val)[:,:0]).hist()
```

```
Out[38]: <matplotlib.axes._subplots.AxesSubplot at 0x7f630bd72a58>
```



## Introduction

This kernel is an attempt to use every trick in the books to unleash the full power of Linear Regression, including a lot of preprocessing and a look at several Regularization algorithms.

At the time of writing, it achieves a score of about 0.121 on the public LB, just using regression, no RF, no xgboost, no ensembling etc. All comments/corrections are more than welcome.

```
In [1]: # Imports
import pandas as pd
import numpy as np
from sklearn.model_selection import cross_val_score, train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LinearRegression, RidgeCV, LassoCV, Elastic
NetCV
from sklearn.metrics import mean_squared_error, make_scorer
from scipy.stats import skew
from IPython.display import display
import matplotlib.pyplot as plt
import seaborn as sns

# Definitions
pd.set_option('display.float_format', lambda x: '%.3f' % x)
%matplotlib inline
#njobs = 4
```

```
In [2]: # Get data
train = pd.read_csv("../input/train.csv")
print("train : " + str(train.shape))

train : (1460, 81)
```

```
In [3]: # Check for duplicates
idsUnique = len(set(train.Id))
idsTotal = train.shape[0]
idsDupli = idsTotal - idsUnique
print("There are " + str(idsDupli) + " duplicate IDs for " + str(idsTotal) +
" total entries")

# Drop Id column
train.drop("Id", axis = 1, inplace = True)

There are 0 duplicate IDs for 1460 total entries
```

## Preprocessing

```
In [4]: # Looking for outliers, as indicated in https://ww2.amstat.org/publications/jse/v19n3/decock.pdf
plt.scatter(train.GrLivArea, train.SalePrice, c = "blue", marker = "s")
plt.title("Looking for outliers")
plt.xlabel("GrLivArea")
plt.ylabel("SalePrice")
plt.show()

train = train[train.GrLivArea < 4000]
```



There seems to be 2 extreme outliers on the bottom right, really large houses that sold for really cheap. More generally, the author of the dataset recommends removing 'any houses with more than 4000 square feet' from the dataset.

Reference : <https://ww2.amstat.org/publications/jse/v19n3/decock.pdf> (<https://ww2.amstat.org/publications/jse/v19n3/decock.pdf>)

```
In [5]: # Log transform the target for official scoring
train.SalePrice = np.log1p(train.SalePrice)
y = train.SalePrice
```

Taking logs means that errors in predicting expensive houses and cheap houses will affect the result equally.

```

In [6]: # Handle missing values for features where median/mean or most common value
        # doesn't make sense

        # Alley : data description says NA means "no alley access"
        train.loc[:, "Alley"] = train.loc[:, "Alley"].fillna("None")
        # BedroomAbvGr : NA most likely means 0
        train.loc[:, "BedroomAbvGr"] = train.loc[:, "BedroomAbvGr"].fillna(0)
        # BsmtQual etc : data description says NA for basement features is "no basem
        # ent"
        train.loc[:, "BsmtQual"] = train.loc[:, "BsmtQual"].fillna("No")
        train.loc[:, "BsmtCond"] = train.loc[:, "BsmtCond"].fillna("No")
        train.loc[:, "BsmtExposure"] = train.loc[:, "BsmtExposure"].fillna("No")
        train.loc[:, "BsmtFinType1"] = train.loc[:, "BsmtFinType1"].fillna("No")
        train.loc[:, "BsmtFinType2"] = train.loc[:, "BsmtFinType2"].fillna("No")
        train.loc[:, "BsmtFullBath"] = train.loc[:, "BsmtFullBath"].fillna(0)
        train.loc[:, "BsmtHalfBath"] = train.loc[:, "BsmtHalfBath"].fillna(0)
        train.loc[:, "BsmtUnfSF"] = train.loc[:, "BsmtUnfSF"].fillna(0)
        # CentralAir : NA most likely means No
        train.loc[:, "CentralAir"] = train.loc[:, "CentralAir"].fillna("N")
        # Condition : NA most likely means Normal
        train.loc[:, "Condition1"] = train.loc[:, "Condition1"].fillna("Norm")
        train.loc[:, "Condition2"] = train.loc[:, "Condition2"].fillna("Norm")
        # EnclosedPorch : NA most likely means no enclosed porch
        train.loc[:, "EnclosedPorch"] = train.loc[:, "EnclosedPorch"].fillna(0)
        # External stuff : NA most likely means average
        train.loc[:, "ExterCond"] = train.loc[:, "ExterCond"].fillna("TA")
        train.loc[:, "ExterQual"] = train.loc[:, "ExterQual"].fillna("TA")
        # Fence : data description says NA means "no fence"
        train.loc[:, "Fence"] = train.loc[:, "Fence"].fillna("No")
        # FireplaceQu : data description says NA means "no fireplace"
        train.loc[:, "FireplaceQu"] = train.loc[:, "FireplaceQu"].fillna("No")
        train.loc[:, "Fireplaces"] = train.loc[:, "Fireplaces"].fillna(0)
        # Functional : data description says NA means typical
        train.loc[:, "Functional"] = train.loc[:, "Functional"].fillna("Typ")
        # GarageType etc : data description says NA for garage features is "no garag
        # e"
        train.loc[:, "GarageType"] = train.loc[:, "GarageType"].fillna("No")
        train.loc[:, "GarageFinish"] = train.loc[:, "GarageFinish"].fillna("No")
        train.loc[:, "GarageQual"] = train.loc[:, "GarageQual"].fillna("No")
        train.loc[:, "GarageCond"] = train.loc[:, "GarageCond"].fillna("No")
        train.loc[:, "GarageArea"] = train.loc[:, "GarageArea"].fillna(0)
        train.loc[:, "GarageCars"] = train.loc[:, "GarageCars"].fillna(0)
        # HalfBath : NA most likely means no half baths above grade
        train.loc[:, "HalfBath"] = train.loc[:, "HalfBath"].fillna(0)
        # HeatingQC : NA most likely means typical
        train.loc[:, "HeatingQC"] = train.loc[:, "HeatingQC"].fillna("TA")
        # KitchenAbvGr : NA most likely means 0
        train.loc[:, "KitchenAbvGr"] = train.loc[:, "KitchenAbvGr"].fillna(0)
        # KitchenQual : NA most likely means typical
        train.loc[:, "KitchenQual"] = train.loc[:, "KitchenQual"].fillna("TA")
        # LotFrontage : NA most likely means no lot frontage
        train.loc[:, "LotFrontage"] = train.loc[:, "LotFrontage"].fillna(0)
        # LotShape : NA most likely means regular
        train.loc[:, "LotShape"] = train.loc[:, "LotShape"].fillna("Reg")
        # MasVnrType : NA most likely means no veneer
        train.loc[:, "MasVnrType"] = train.loc[:, "MasVnrType"].fillna("None")
        train.loc[:, "MasVnrArea"] = train.loc[:, "MasVnrArea"].fillna(0)
        # MiscFeature : data description says NA means "no misc feature"
        train.loc[:, "MiscFeature"] = train.loc[:, "MiscFeature"].fillna("No")
        train.loc[:, "MiscVal"] = train.loc[:, "MiscVal"].fillna(0)
        # OpenPorchSF : NA most likely means no open porch
        train.loc[:, "OpenPorchSF"] = train.loc[:, "OpenPorchSF"].fillna(0)
        # PavedDrive : NA most likely means not paved
        train.loc[:, "PavedDrive"] = train.loc[:, "PavedDrive"].fillna("N")
        # PoolQC : data description says NA means "no pool"
        train.loc[:, "PoolQC"] = train.loc[:, "PoolQC"].fillna("No")
        train.loc[:, "PoolArea"] = train.loc[:, "PoolArea"].fillna(0)
        # SaleCondition : NA most likely means normal sale

```

```
In [7]: # Some numerical features are actually really categories
train = train.replace({"MSSubClass" : {20 : "SC20", 30 : "SC30", 40 : "SC40", 45 : "SC45",
                                         50 : "SC50", 60 : "SC60", 70 : "SC70", 75 : "SC75",
                                         80 : "SC80", 85 : "SC85", 90 : "SC90", 120 : "SC120",
                                         150 : "SC150", 160 : "SC160", 180 : "SC180", 190 : "SC190"},
                      "MoSold" : {1 : "Jan", 2 : "Feb", 3 : "Mar", 4 : "Apr", 5 : "May", 6 : "Jun",
                                   7 : "Jul", 8 : "Aug", 9 : "Sep", 10 : "Oct", 11 : "Nov", 12 : "Dec"}}
                      })
```

```
In [8]: # Encode some categorical features as ordered numbers when there is information in the order
train = train.replace({"Alley" : {"Grvl" : 1, "Pave" : 2},
                      "BsmtCond" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "BsmtExposure" : {"No" : 0, "Mn" : 1, "Av" : 2, "Gd" : 3},
                      "BsmtFinType1" : {"No" : 0, "Unf" : 1, "LwQ" : 2, "Re" : 3, "BLQ" : 4,
                                           "ALQ" : 5, "GLQ" : 6},
                      "BsmtFinType2" : {"No" : 0, "Unf" : 1, "LwQ" : 2, "Re" : 3, "BLQ" : 4,
                                           "ALQ" : 5, "GLQ" : 6},
                      "BsmtQual" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "ExterCond" : {"Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "ExterQual" : {"Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "FireplaceQu" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "Functional" : {"Sal" : 1, "Sev" : 2, "Maj2" : 3, "Maj1" : 4, "Mod" : 5,
                                       "Min2" : 6, "Min1" : 7, "Typ" : 8},
                      "GarageCond" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "GarageQual" : {"No" : 0, "Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "HeatingQC" : {"Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "KitchenQual" : {"Po" : 1, "Fa" : 2, "TA" : 3, "Gd" : 4, "Ex" : 5},
                      "LandSlope" : {"Sev" : 1, "Mod" : 2, "Gtl" : 3},
                      "LotShape" : {"IR3" : 1, "IR2" : 2, "IR1" : 3, "Reg" : 4},
                      "PavedDrive" : {"N" : 0, "P" : 1, "Y" : 2},
                      "PoolQC" : {"No" : 0, "Fa" : 1, "TA" : 2, "Gd" : 3, "Ex" : 4},
                      "Street" : {"Grvl" : 1, "Pave" : 2},
                      "Utilities" : {"ELO" : 1, "NoSeWa" : 2, "NoSewr" : 3, "AllPub" : 4}}
                      )
```

Then we will create new features, in 3 ways :

1. Simplifications of existing features
2. Combinations of existing features
3. Polynomials on the top 10 existing features

```

In [9]: # Create new features
# 1* Simplifications of existing features
train["SimplOverallQual"] = train.OverallQual.replace({1 : 1, 2 : 1, 3 : 1,
# bad
                                                    4 : 2, 5 : 2, 6 : 2,
# average
                                                    7 : 3, 8 : 3, 9 : 3,
10 : 3 # good
                                                    })
train["SimplOverallCond"] = train.OverallCond.replace({1 : 1, 2 : 1, 3 : 1,
# bad
                                                    4 : 2, 5 : 2, 6 : 2,
# average
                                                    7 : 3, 8 : 3, 9 : 3,
10 : 3 # good
                                                    })
train["SimplPoolQC"] = train.PoolQC.replace({1 : 1, 2 : 1, # average
3 : 2, 4 : 2 # good
                                                    })
train["SimplGarageCond"] = train.GarageCond.replace({1 : 1, # bad
2 : 1, 3 : 1, # average
4 : 2, 5 : 2 # good
                                                    })
train["SimplGarageQual"] = train.GarageQual.replace({1 : 1, # bad
2 : 1, 3 : 1, # average
4 : 2, 5 : 2 # good
                                                    })
train["SimplFireplaceQu"] = train.FireplaceQu.replace({1 : 1, # bad
2 : 1, 3 : 1, # average
4 : 2, 5 : 2 # good
                                                    })
train["SimplFireplaceQu"] = train.FireplaceQu.replace({1 : 1, # bad
2 : 1, 3 : 1, # average
4 : 2, 5 : 2 # good
                                                    })
train["SimplFunctional"] = train.Functional.replace({1 : 1, 2 : 1, # bad
3 : 2, 4 : 2, # major
5 : 3, 6 : 3, 7 : 3, #
minor
8 : 4 # typical
                                                    })
train["SimplKitchenQual"] = train.KitchenQual.replace({1 : 1, # bad
2 : 1, 3 : 1, # average
4 : 2, 5 : 2 # good
                                                    })
train["SimplHeatingQC"] = train.HeatingQC.replace({1 : 1, # bad
2 : 1, 3 : 1, # average
4 : 2, 5 : 2 # good
                                                    })
train["SimplBsmtFinType1"] = train.BsmtFinType1.replace({1 : 1, # unfinished
2 : 1, 3 : 1, # rec
room
4 : 2, 5 : 2, 6 : 2
# living quarters
                                                    })
train["SimplBsmtFinType2"] = train.BsmtFinType2.replace({1 : 1, # unfinished
2 : 1, 3 : 1, # rec
room
4 : 2, 5 : 2, 6 : 2
# living quarters
                                                    })
train["SimplBsmtCond"] = train.BsmtCond.replace({1 : 1, # bad
2 : 1, 3 : 1, # average
4 : 2, 5 : 2 # good
                                                    })

```

```
In [10]: # Find most important features relative to target
print("Find most important features relative to target")
corr = train.corr()
corr.sort_values(["SalePrice"], ascending = False, inplace = True)
print(corr.SalePrice)
```



Find most important features relative to target

SalePrice	1.000
OverallQual	0.819
AllSF	0.817
AllFlrsSF	0.729
GrLivArea	0.719
SimplOverallQual	0.708
ExterQual	0.681
GarageCars	0.680
TotalBath	0.673
KitchenQual	0.667
GarageScore	0.657
GarageArea	0.655
TotalBsmtSF	0.642
SimplExterQual	0.636
SimplGarageScore	0.631
BsmtQual	0.615
1stFlrSF	0.614
SimplKitchenQual	0.610
OverallGrade	0.604
SimplBsmtQual	0.594
FullBath	0.591
YearBuilt	0.589
ExterGrade	0.587
YearRemodAdd	0.569
FireplaceQu	0.547
GarageYrBlt	0.544
TotRmsAbvGrd	0.533
SimplOverallGrade	0.527
SimplKitchenScore	0.523
FireplaceScore	0.518
...	
SimplBsmtCond	0.204
BedroomAbvGr	0.204
AllPorchSF	0.199
LotFrontage	0.174
SimplFunctional	0.137
Functional	0.136
ScreenPorch	0.124
SimplBsmtFinType2	0.105
Street	0.058
3SsnPorch	0.056
ExterCond	0.051
PoolArea	0.041
SimplPoolScore	0.040
SimplPoolQC	0.040
PoolScore	0.040
PoolQC	0.038
BsmtFinType2	0.016
Utilities	0.013
BsmtFinSF2	0.006
BsmtHalfBath	-0.015
MiscVal	-0.020
SimplOverallCond	-0.028
YrSold	-0.034
OverallCond	-0.037
LowQualFinSF	-0.038
LandSlope	-0.040
SimplExterCond	-0.042
KitchenAbvGr	-0.148
EnclosedPorch	-0.149
LotShape	-0.286

Name: SalePrice, dtype: float64

```
In [11]: # Create new features
# 3* Polynomials on the top 10 existing features
train["OverallQual-s2"] = train["OverallQual"] ** 2
train["OverallQual-s3"] = train["OverallQual"] ** 3
train["OverallQual-Sq"] = np.sqrt(train["OverallQual"])
train["AllSF-2"] = train["AllSF"] ** 2
train["AllSF-3"] = train["AllSF"] ** 3
train["AllSF-Sq"] = np.sqrt(train["AllSF"])
train["AllFlrsSF-2"] = train["AllFlrsSF"] ** 2
train["AllFlrsSF-3"] = train["AllFlrsSF"] ** 3
train["AllFlrsSF-Sq"] = np.sqrt(train["AllFlrsSF"])
train["GrLivArea-2"] = train["GrLivArea"] ** 2
train["GrLivArea-3"] = train["GrLivArea"] ** 3
train["GrLivArea-Sq"] = np.sqrt(train["GrLivArea"])
train["SimplOverallQual-s2"] = train["SimplOverallQual"] ** 2
train["SimplOverallQual-s3"] = train["SimplOverallQual"] ** 3
train["SimplOverallQual-Sq"] = np.sqrt(train["SimplOverallQual"])
train["ExterQual-2"] = train["ExterQual"] ** 2
train["ExterQual-3"] = train["ExterQual"] ** 3
train["ExterQual-Sq"] = np.sqrt(train["ExterQual"])
train["GarageCars-2"] = train["GarageCars"] ** 2
train["GarageCars-3"] = train["GarageCars"] ** 3
train["GarageCars-Sq"] = np.sqrt(train["GarageCars"])
train["TotalBath-2"] = train["TotalBath"] ** 2
train["TotalBath-3"] = train["TotalBath"] ** 3
train["TotalBath-Sq"] = np.sqrt(train["TotalBath"])
train["KitchenQual-2"] = train["KitchenQual"] ** 2
train["KitchenQual-3"] = train["KitchenQual"] ** 3
train["KitchenQual-Sq"] = np.sqrt(train["KitchenQual"])
train["GarageScore-2"] = train["GarageScore"] ** 2
train["GarageScore-3"] = train["GarageScore"] ** 3
train["GarageScore-Sq"] = np.sqrt(train["GarageScore"])
```

```
In [12]: # Differentiate numerical features (minus the target) and categorical features
categorical_features = train.select_dtypes(include = ["object"]).columns
numerical_features = train.select_dtypes(exclude = ["object"]).columns
numerical_features = numerical_features.drop("SalePrice")
print("Numerical features : " + str(len(numerical_features)))
print("Categorical features : " + str(len(categorical_features)))
train_num = train[numerical_features]
train_cat = train[categorical_features]
```

Numerical features : 117  
Categorical features : 26

```
In [13]: # Handle remaining missing values for numerical features by using median as replacement
print("NAs for numerical features in train : " + str(train_num.isnull().values.sum()))
train_num = train_num.fillna(train_num.median())
print("Remaining NAs for numerical features in train : " + str(train_num.isnull().values.sum()))
```

NAs for numerical features in train : 81  
Remaining NAs for numerical features in train : 0

```
In [14]: # Log transform of the skewed numerical features to lessen impact of outliers
# Inspired by Alexandru Papiu's script : https://www.kaggle.com/apapiu/house-prices-advanced-regression-techniques/regularized-linear-models
# As a general rule of thumb, a skewness with an absolute value > 0.5 is considered at least moderately skewed
skewness = train_num.apply(Lambda x: skew(x))
skewness = skewness[abs(skewness) > 0.5]
print(str(skewness.shape[0]) + " skewed numerical features to log transform")
skewed_features = skewness.index
train_num[skewed_features] = np.log1p(train_num[skewed_features])
```

86 skewed numerical features to log transform

```
In [15]: # Create dummy features for categorical values via one-hot encoding
print("NAs for categorical features in train : " + str(train_cat.isnull().values.sum()))
train_cat = pd.get_dummies(train_cat)
print("Remaining NAs for categorical features in train : " + str(train_cat.isnull().values.sum()))
```

NAs for categorical features in train : 1

Remaining NAs for categorical features in train : 0

## Modeling

```
In [16]: # Join categorical and numerical features
train = pd.concat([train_num, train_cat], axis = 1)
print("New number of features : " + str(train.shape[1]))

# Partition the dataset in train + validation sets
X_train, X_test, y_train, y_test = train_test_split(train, y, test_size = 0.3, random_state = 0)
print("X_train : " + str(X_train.shape))
print("X_test : " + str(X_test.shape))
print("y_train : " + str(y_train.shape))
print("y_test : " + str(y_test.shape))
```

New number of features : 319

X\_train : (1019, 319)

X\_test : (437, 319)

y\_train : (1019,)

y\_test : (437,)

```
In [17]: # Standardize numerical features
stdSc = StandardScaler()
X_train.loc[:, numerical_features] = stdSc.fit_transform(X_train.loc[:, numerical_features])
X_test.loc[:, numerical_features] = stdSc.transform(X_test.loc[:, numerical_features])
```

/opt/conda/lib/python3.5/site-packages/pandas/core/indexing.py:465: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
self.obj[item] = s
```

Standardization cannot be done before the partitioning, as we don't want to fit the StandardScaler on some observations that will later be used in the test set.

```
In [18]: # Define error measure for official scoring : RMSE
scorer = make_scorer(mean_squared_error, greater_is_better = False)

def rmse_cv_train(model):
    rmse= np.sqrt(-cross_val_score(model, X_train, y_train, scoring = scorer,
    cv = 10))
    return(rmse)

def rmse_cv_test(model):
    rmse= np.sqrt(-cross_val_score(model, X_test, y_test, scoring = scorer,
    cv = 10))
    return(rmse)
```

Note : I'm not getting nearly the same numbers in local CV compared to public LB, so I'm a tad worried that my CV process may have an issue somewhere. If you spot something, please let me know.

### 1\* Linear Regression without regularization

```
In [19]: # Linear Regression
lr = LinearRegression()
lr.fit(X_train, y_train)

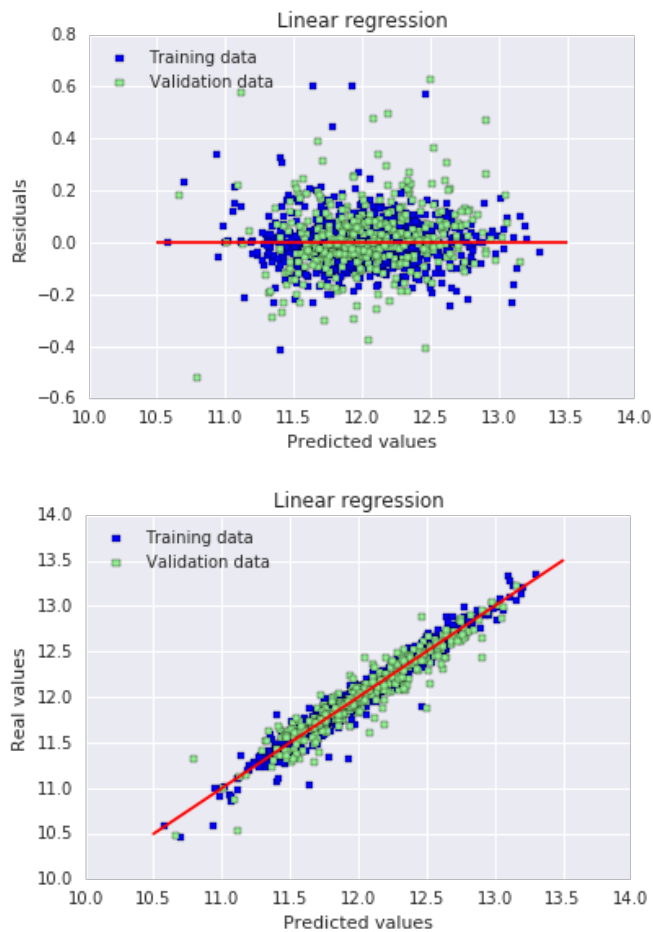
# Look at predictions on training and validation set
print("RMSE on Training set :", rmse_cv_train(lr).mean())
print("RMSE on Test set :", rmse_cv_test(lr).mean())
y_train_pred = lr.predict(X_train)
y_test_pred = lr.predict(X_test)

# Plot residuals
plt.scatter(y_train_pred, y_train_pred - y_train, c = "blue", marker = "s",
label = "Training data")
plt.scatter(y_test_pred, y_test_pred - y_test, c = "lightgreen", marker = "s",
label = "Validation data")
plt.title("Linear regression")
plt.xlabel("Predicted values")
plt.ylabel("Residuals")
plt.legend(loc = "upper left")
plt.hlines(y = 0, xmin = 10.5, xmax = 13.5, color = "red")
plt.show()

# Plot predictions
plt.scatter(y_train_pred, y_train, c = "blue", marker = "s", label = "Training data")
plt.scatter(y_test_pred, y_test, c = "lightgreen", marker = "s", label = "Validation data")
plt.title("Linear regression")
plt.xlabel("Predicted values")
plt.ylabel("Real values")
plt.legend(loc = "upper left")
plt.plot([10.5, 13.5], [10.5, 13.5], c = "red")
plt.show()
```

RMSE on Training set : 15758371373.7

RMSE on Test set : 0.395779797728



RMSE on Training set shows up weird here (not when I run it on my computer) for some reason.

Errors seem randomly distributed and randomly scattered around the centerline, so there is that at least. It means our model was able to capture most of the explanatory information.

## 2\* Linear Regression with Ridge regularization (L2 penalty)

From the *Python Machine Learning* book by Sebastian Raschka : Regularization is a very useful method to handle collinearity, filter out noise from data, and eventually prevent overfitting. The concept behind regularization is to introduce additional information (bias) to penalize extreme parameter weights.

Ridge regression is an L2 penalized model where we simply add the squared sum of the weights to our cost function.

```

In [20]: # 2* Ridge
ridge = RidgeCV(alphas = [0.01, 0.03, 0.06, 0.1, 0.3, 0.6, 1, 3, 6, 10, 30,
60])
ridge.fit(X_train, y_train)
alpha = ridge.alpha_
print("Best alpha :", alpha)

print("Try again for more precision with alphas centered around " + str(alpha))
ridge = RidgeCV(alphas = [alpha * .6, alpha * .65, alpha * .7, alpha * .75,
alpha * .8, alpha * .85,
alpha * .9, alpha * .95, alpha, alpha * 1.05, alpha * 1.1, alpha * 1.15,
alpha * 1.25, alpha * 1.3, alpha * 1.35, alpha * 1.4],
cv = 10)
ridge.fit(X_train, y_train)
alpha = ridge.alpha_
print("Best alpha :", alpha)

print("Ridge RMSE on Training set :", rmse_cv_train(ridge).mean())
print("Ridge RMSE on Test set :", rmse_cv_test(ridge).mean())
y_train_rdg = ridge.predict(X_train)
y_test_rdg = ridge.predict(X_test)

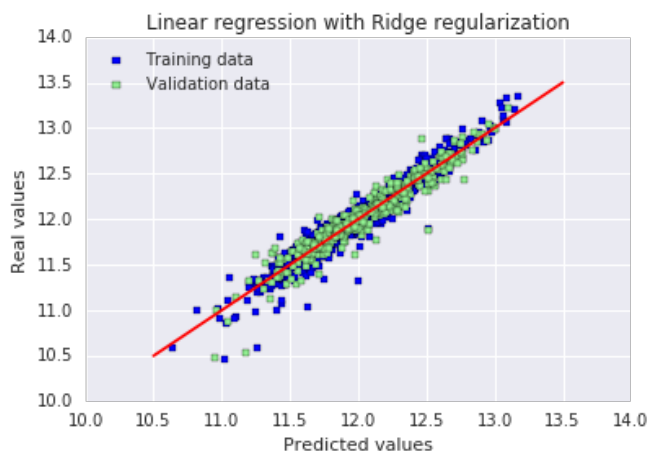
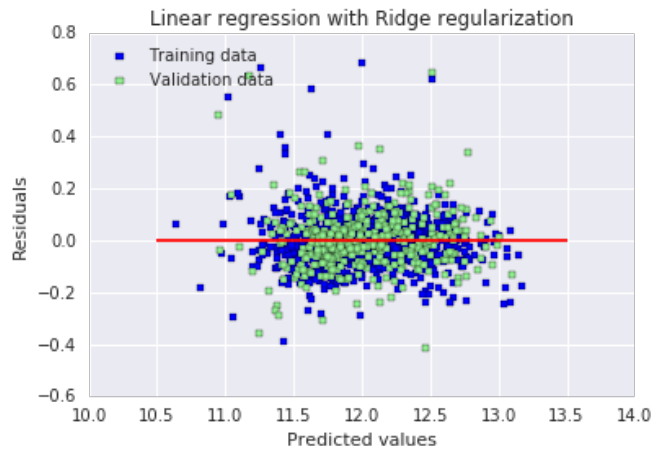
# Plot residuals
plt.scatter(y_train_rdg, y_train_rdg - y_train, c = "blue", marker = "s", label = "Training data")
plt.scatter(y_test_rdg, y_test_rdg - y_test, c = "lightgreen", marker = "s", label = "Validation data")
plt.title("Linear regression with Ridge regularization")
plt.xlabel("Predicted values")
plt.ylabel("Residuals")
plt.legend(loc = "upper left")
plt.hlines(y = 0, xmin = 10.5, xmax = 13.5, color = "red")
plt.show()

# Plot predictions
plt.scatter(y_train_rdg, y_train, c = "blue", marker = "s", label = "Training data")
plt.scatter(y_test_rdg, y_test, c = "lightgreen", marker = "s", label = "Validation data")
plt.title("Linear regression with Ridge regularization")
plt.xlabel("Predicted values")
plt.ylabel("Real values")
plt.legend(loc = "upper left")
plt.plot([10.5, 13.5], [10.5, 13.5], c = "red")
plt.show()

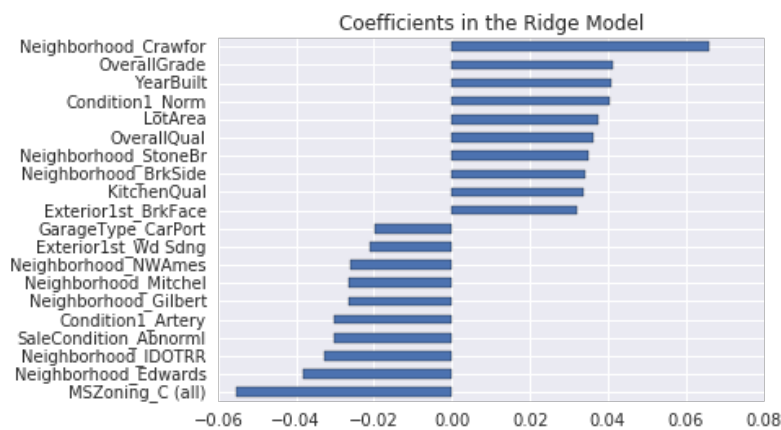
# Plot important coefficients
coefs = pd.Series(ridge.coef_, index = X_train.columns)
print("Ridge picked " + str(sum(coefs != 0)) + " features and eliminated the other " + \
str(sum(coefs == 0)) + " features")
imp_coefs = pd.concat([coefs.sort_values().head(10),
coefs.sort_values().tail(10)])
imp_coefs.plot(kind = "barh")
plt.title("Coefficients in the Ridge Model")
plt.show()

```

Best alpha : 30.0  
 Try again for more precision with alphas centered around 30.0  
 Best alpha : 24.0  
 Ridge RMSE on Training set : 0.115405723285  
 Ridge RMSE on Test set : 0.116427213778



Ridge picked 316 features and eliminated the other 3 features



We're getting a much better RMSE result now that we've added regularization. The very small difference between training and test results indicate that we eliminated most of the overfitting. Visually, the graphs seem to confirm that idea.

Ridge used almost all of the existing features.



**3\* Linear Regression with Lasso regularization (L1 penalty)**

LASSO stands for *Least Absolute Shrinkage and Selection Operator*. It is an alternative regularization method, where we simply replace the square of the weights by the sum of the absolute value of the weights. In contrast to L2 regularization, L1 regularization yields sparse feature vectors : most feature weights will be zero. Sparsity can be useful in practice if we have a high dimensional dataset with many features that are irrelevant.

We can suspect that it should be more efficient than Ridge here.

```

In [21]: # 3* Lasso
lasso = LassoCV(alphas = [0.0001, 0.0003, 0.0006, 0.001, 0.003, 0.006, 0.01,
                          0.03, 0.06, 0.1,
                          0.3, 0.6, 1],
                max_iter = 50000, cv = 10)
lasso.fit(X_train, y_train)
alpha = lasso.alpha_
print("Best alpha :", alpha)

print("Try again for more precision with alphas centered around " + str(alpha))
lasso = LassoCV(alphas = [alpha * .6, alpha * .65, alpha * .7, alpha * .75,
                          alpha * .8,
                          alpha * .85, alpha * .9, alpha * .95, alpha, alpha
                          * 1.05,
                          alpha * 1.1, alpha * 1.15, alpha * 1.25, alpha *
                          1.3, alpha * 1.35,
                          alpha * 1.4],
                max_iter = 50000, cv = 10)
lasso.fit(X_train, y_train)
alpha = lasso.alpha_
print("Best alpha :", alpha)

print("Lasso RMSE on Training set :", rmse_cv_train(lasso).mean())
print("Lasso RMSE on Test set :", rmse_cv_test(lasso).mean())
y_train_las = lasso.predict(X_train)
y_test_las = lasso.predict(X_test)

# Plot residuals
plt.scatter(y_train_las, y_train_las - y_train, c = "blue", marker = "s", label = "Training data")
plt.scatter(y_test_las, y_test_las - y_test, c = "lightgreen", marker = "s", label = "Validation data")
plt.title("Linear regression with Lasso regularization")
plt.xlabel("Predicted values")
plt.ylabel("Residuals")
plt.legend(loc = "upper left")
plt.hlines(y = 0, xmin = 10.5, xmax = 13.5, color = "red")
plt.show()

# Plot predictions
plt.scatter(y_train_las, y_train, c = "blue", marker = "s", label = "Training data")
plt.scatter(y_test_las, y_test, c = "lightgreen", marker = "s", label = "Validation data")
plt.title("Linear regression with Lasso regularization")
plt.xlabel("Predicted values")
plt.ylabel("Real values")
plt.legend(loc = "upper left")
plt.plot([10.5, 13.5], [10.5, 13.5], c = "red")
plt.show()

# Plot important coefficients
coefs = pd.Series(lasso.coef_, index = X_train.columns)
print("Lasso picked " + str(sum(coefs != 0)) + " features and eliminated the other " + \
      str(sum(coefs == 0)) + " features")
imp_coefs = pd.concat([coefs.sort_values().head(10),
                      coefs.sort_values().tail(10)])
imp_coefs.plot(kind = "barh")
plt.title("Coefficients in the Lasso Model")
plt.show()

```

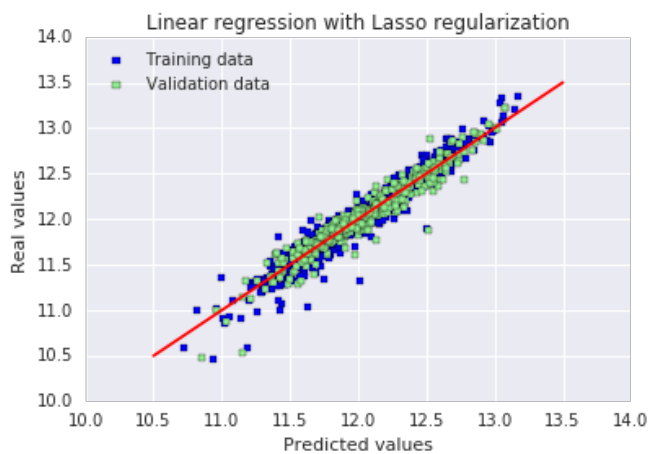
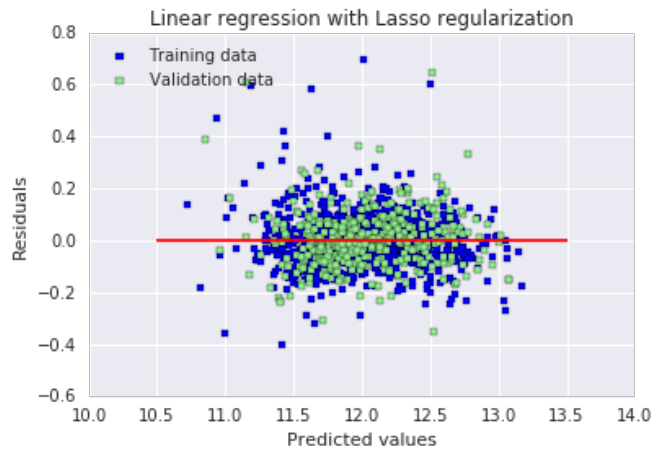
Best alpha : 0.0006

Try again for more precision with alphas centered around 0.0006

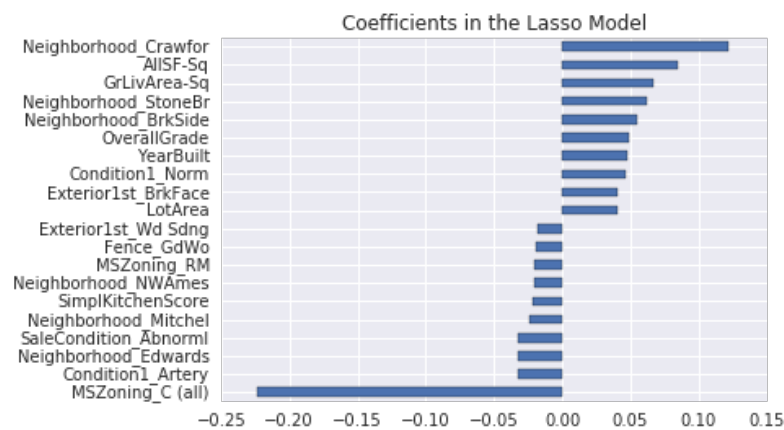
Best alpha : 0.0006

Lasso RMSE on Training set : 0.114111508375

Lasso RMSE on Test set : 0.115832132218



Lasso picked 110 features and eliminated the other 209 features



RMSE results are better both on training and test sets. The most interesting thing is that Lasso used only one third of the available features. Another interesting tidbit : it seems to give big weights to Neighborhood categories, both in positive and negative ways. Intuitively it makes sense, house prices change a whole lot from one neighborhood to another in the same city.

The "MSZoning\_C (all)" feature seems to have a disproportionate impact compared to the others. It is defined as *general zoning classification : commercial*. It seems a bit weird to me that having your house in a mostly commercial zone would be such a terrible thing.

#### **4\* Linear Regression with ElasticNet regularization (L1 and L2 penalty)**

ElasticNet is a compromise between Ridge and Lasso regression. It has a L1 penalty to generate sparsity and a L2 penalty to overcome some of the limitations of Lasso, such as the number of variables (Lasso can't select more features than it has observations, but it's not the case here anyway).

```

In [22]: # 4* ElasticNet
elasticNet = ElasticNetCV(l1_ratio = [0.1, 0.3, 0.5, 0.6, 0.7, 0.8, 0.85, 0.9, 0.95, 1],
                           alphas = [0.0001, 0.0003, 0.0006, 0.001, 0.003, 0.006,
                                       0.01, 0.03, 0.06, 0.1, 0.3, 0.6, 1, 3, 6],
                           max_iter = 50000, cv = 10)
elasticNet.fit(X_train, y_train)
alpha = elasticNet.alpha_
ratio = elasticNet.l1_ratio_
print("Best l1_ratio :", ratio)
print("Best alpha :", alpha )

print("Try again for more precision with l1_ratio centered around " + str(ratio))
elasticNet = ElasticNetCV(l1_ratio = [ratio * .85, ratio * .9, ratio * .95, ratio, ratio * 1.05, ratio * 1.1, ratio * 1.15],
                           alphas = [0.0001, 0.0003, 0.0006, 0.001, 0.003, 0.006, 0.01, 0.03, 0.06, 0.1, 0.3, 0.6, 1, 3, 6],
                           max_iter = 50000, cv = 10)
elasticNet.fit(X_train, y_train)
if (elasticNet.l1_ratio_ > 1):
    elasticNet.l1_ratio_ = 1
alpha = elasticNet.alpha_
ratio = elasticNet.l1_ratio_
print("Best l1_ratio :", ratio)
print("Best alpha :", alpha )

print("Now try again for more precision on alpha, with l1_ratio fixed at " + str(ratio) +
      " and alpha centered around " + str(alpha))
elasticNet = ElasticNetCV(l1_ratio = ratio,
                           alphas = [alpha * .6, alpha * .65, alpha * .7, alpha * .75, alpha * .8, alpha * .85, alpha * .9, alpha * .95, alpha, alpha * 1.05, alpha * 1.1, alpha * 1.15, alpha * 1.25, alpha * 1.3, alpha * 1.35, alpha * 1.4],
                           max_iter = 50000, cv = 10)
elasticNet.fit(X_train, y_train)
if (elasticNet.l1_ratio_ > 1):
    elasticNet.l1_ratio_ = 1
alpha = elasticNet.alpha_
ratio = elasticNet.l1_ratio_
print("Best l1_ratio :", ratio)
print("Best alpha :", alpha )

print("ElasticNet RMSE on Training set :", rmse_cv_train(elasticNet).mean())
print("ElasticNet RMSE on Test set :", rmse_cv_test(elasticNet).mean())
y_train_ela = elasticNet.predict(X_train)
y_test_ela = elasticNet.predict(X_test)

# Plot residuals
plt.scatter(y_train_ela, y_train_ela - y_train, c = "blue", marker = "s", label = "Training data")
plt.scatter(y_test_ela, y_test_ela - y_test, c = "lightgreen", marker = "s", label = "Validation data")
plt.title("Linear regression with ElasticNet regularization")
plt.xlabel("Predicted values")
plt.ylabel("Residuals")
plt.legend(loc = "upper left")
plt.hlines(y = 0, xmin = 10.5, xmax = 13.5, color = "red")
plt.show()

# Plot predictions
plt.scatter(y_train, y_train_ela, c = "blue", marker = "s", label = "Training data")
plt.scatter(y_test, y_test_ela, c = "lightgreen", marker = "s", label = "Validation data")

```

```
/opt/conda/lib/python3.5/site-packages/sklearn/linear_model/coordinate_descen
t.py:479: ConvergenceWarning: Objective did not converge. You might want to i
ncrease the number of iterations. Fitting data with very small alpha may caus
e precision problems.
  ConvergenceWarning)
```

Best l1\_ratio : 1.0

Best alpha : 0.0006

Try again for more precision with l1\_ratio centered around 1.0

Best l1\_ratio : 1.0

Best alpha : 0.0006

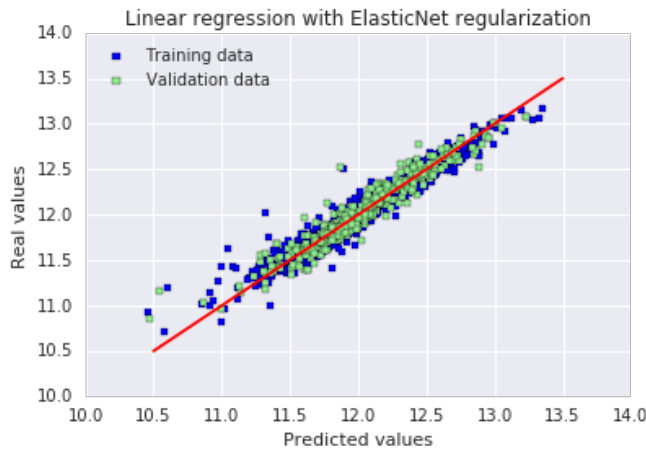
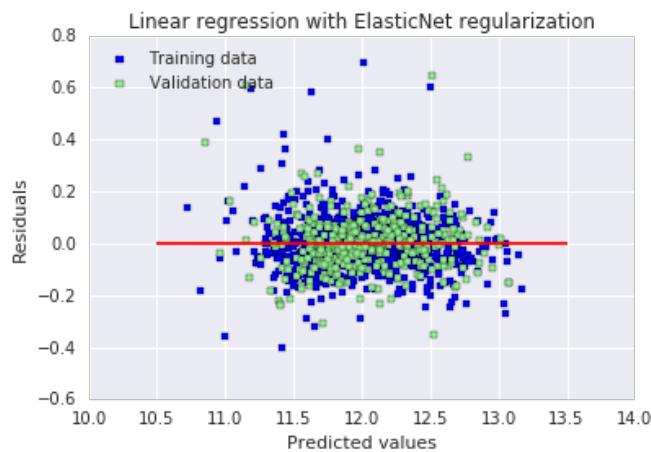
Now try again for more precision on alpha, with l1\_ratio fixed at 1.0 and alp  
ha centered around 0.0006

Best l1\_ratio : 1.0

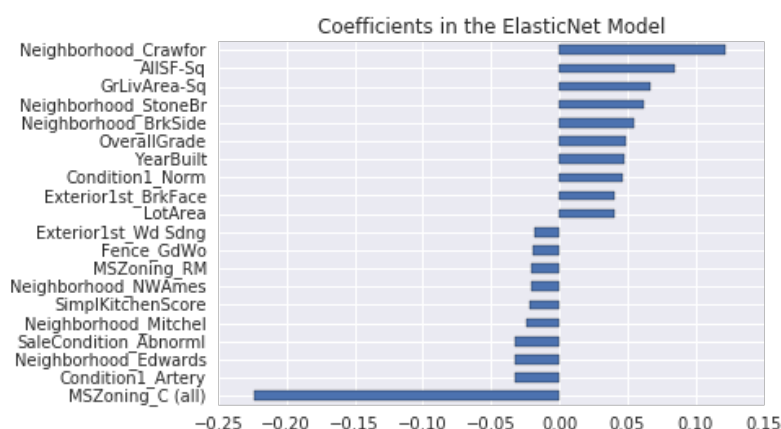
Best alpha : 0.0006

ElasticNet RMSE on Training set : 0.114111508375

ElasticNet RMSE on Test set : 0.115832132218



ElasticNet picked 110 features and eliminated the other 209 features



The optimal L1 ratio used by ElasticNet here is equal to 1, which means it is exactly equal to the Lasso regressor we used earlier (and had it been equal to 0, it would have been exactly equal to our Ridge regressor). The model didn't need any L2 regularization to overcome any potential L1 shortcoming.

Note : I tried to remove the "MSZoning\_C (all)" feature, it resulted in a slightly worse CV score, but slightly better public LB score.

## Conclusion

Putting time and effort into preparing the dataset and optimizing the regularization resulted in a decent score, better than some public scripts which use algorithms that historically perform better in Kaggle contests, like Random Forests. Being fairly new to the world of machine learning contests, I will appreciate any constructive pointer to improve, and I thank you for your time.

# Stacked Regressions to predict House Prices

## Serigne

July 2017

If you use parts of this notebook in your scripts/notebooks, giving some kind of credit would be very much appreciated :) You can for instance link back to this notebook. Thanks!

This competition is very important to me as it helped me to begin my journey on Kaggle few months ago. I've read some great notebooks here. To name a few:

1. [Comprehensive data exploration with Python](https://www.kaggle.com/pmarcelino/comprehensive-data-exploration-with-python) (<https://www.kaggle.com/pmarcelino/comprehensive-data-exploration-with-python>) by **Pedro Marcelino** : Great and very motivational data analysis
2. [A study on Regression applied to the Ames dataset](https://www.kaggle.com/julienca/a-study-on-regression-applied-to-the-ames-dataset) (<https://www.kaggle.com/julienca/a-study-on-regression-applied-to-the-ames-dataset>) by **Julien Cohen-Solal** : Thorough features engineering and deep dive into linear regression analysis but really easy to follow for beginners.
3. [Regularized Linear Models](https://www.kaggle.com/apapiu/regularized-linear-models) (<https://www.kaggle.com/apapiu/regularized-linear-models>) by **Alexandru Papiu** : Great Starter kernel on modelling and Cross-validation

I can't recommend enough every beginner to go carefully through these kernels (and of course through many others great kernels) and get their first insights in data science and kaggle competitions.

After that (and some basic practices) you should be more confident to go through [this great script](https://www.kaggle.com/humananalog/xgboost-lasso) (<https://www.kaggle.com/humananalog/xgboost-lasso>) by **Human Analog** who did an impressive work on features engineering.

As the dataset is particularly handy, I decided few days ago to get back in this competition and apply things I learnt so far, especially stacking models. For that purpose, we build two stacking classes ( the simplest approach and a less simple one).

As these classes are written for general purpose, you can easily adapt them and/or extend them for your regression problems. The overall approach is hopefully concise and easy to follow..

The features engineering is rather parsimonious (at least compared to some others great scripts) . It is pretty much :

- **Imputing missing values** by proceeding sequentially through the data
- **Transforming** some numerical variables that seem really categorical
- **Label Encoding** some categorical variables that may contain information in their ordering set
- **Box Cox Transformation** (<http://onlinestatbook.com/2/transformations/box-cox.html>) of skewed features (instead of log-transformation) : This gave me a **slightly better result** both on leaderboard and cross-validation.
- **Getting dummy variables** for categorical features.

Then we choose many base models (mostly sklearn based models + sklearn API of DMLC's [XGBoost](https://github.com/dmlc/xgboost) (<https://github.com/dmlc/xgboost>) and Microsoft's [LightGBM](https://github.com/Microsoft/LightGBM) (<https://github.com/Microsoft/LightGBM>)), cross-validate them on the data before stacking/ensembling them. The key here is to make the (linear) models robust to outliers. This improved the result both on LB and cross-validation.

To my surprise, this does well on LB ( 0.11420 and top 4% the last time I tested it : **July 2, 2017** )

**Hope that at the end of this notebook, stacking will be clear for those, like myself, who found the concept not so easy to grasp**



```
In [1]: #import some necessary librairies

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
%matplotlib inline
import matplotlib.pyplot as plt # Matlab-style plotting
import seaborn as sns
color = sns.color_palette()
sns.set_style('darkgrid')
import warnings
def ignore_warn(*args, **kwargs):
    pass
warnings.warn = ignore_warn #ignore annoying warning (from sklearn and seaborn)

from scipy import stats
from scipy.stats import norm, skew #for some statistics

pd.set_option('display.float_format', lambda x: '{:.3f}'.format(x)) #Limiting floats output to 3 decimal points

from subprocess import check_output
print(check_output(["ls", "../input"]).decode("utf8")) #check the files available in the directory

sample_submission.csv
test.csv
train.csv
```

```
In [2]: #Now let's import and put the train and test datasets in pandas dataframe

train = pd.read_csv('../input/train.csv')
test = pd.read_csv('../input/test.csv')
```

```
In [3]: ##display the first five rows of the train dataset.
train.head(5)
```

Out[3]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	...
0	1	60	RL	65.000	8450	Pave	NaN	Reg	Lvl	AllPub	...
1	2	20	RL	80.000	9600	Pave	NaN	Reg	Lvl	AllPub	...
2	3	60	RL	68.000	11250	Pave	NaN	IR1	Lvl	AllPub	...
3	4	70	RL	60.000	9550	Pave	NaN	IR1	Lvl	AllPub	...
4	5	60	RL	84.000	14260	Pave	NaN	IR1	Lvl	AllPub	...

5 rows × 81 columns

```
In [4]: ##display the first five rows of the test dataset.
test.head(5)
```

Out[4]:

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	.
0	1461	20	RH	80.000	11622	Pave	NaN	Reg	Lvl	AllPub	.
1	1462	20	RL	81.000	14267	Pave	NaN	IR1	Lvl	AllPub	.
2	1463	60	RL	74.000	13830	Pave	NaN	IR1	Lvl	AllPub	.
3	1464	60	RL	78.000	9978	Pave	NaN	IR1	Lvl	AllPub	.
4	1465	120	RL	43.000	5005	Pave	NaN	IR1	HLS	AllPub	.

5 rows × 80 columns

```
In [5]: #check the numbers of samples and features
print("The train data size before dropping Id feature is : {}".format(train.shape))
print("The test data size before dropping Id feature is : {}".format(test.shape))

#Save the 'Id' column
train_ID = train['Id']
test_ID = test['Id']

#Now drop the 'Id' colum since it's unnecessary for the prediction process.
train.drop("Id", axis = 1, inplace = True)
test.drop("Id", axis = 1, inplace = True)

#check again the data size after dropping the 'Id' variable
print("\nThe train data size after dropping Id feature is : {}".format(train.shape))
print("The test data size after dropping Id feature is : {}".format(test.shape))
```

The train data size before dropping Id feature is : (1460, 81)  
The test data size before dropping Id feature is : (1459, 80)

The train data size after dropping Id feature is : (1460, 80)  
The test data size after dropping Id feature is : (1459, 79)

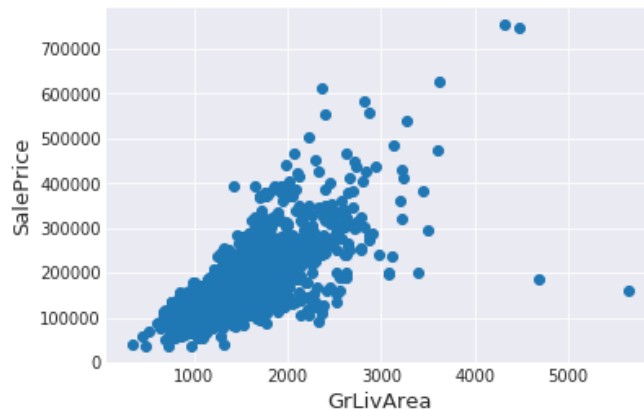
## Data Processing

### Outliers

[Documentation \(http://ww2.amstat.org/publications/jse/v19n3/Decock/DataDocumentation.txt\)](http://ww2.amstat.org/publications/jse/v19n3/Decock/DataDocumentation.txt) for the Ames Housing Data indicates that there are outliers present in the training data

Let's explore these outliers

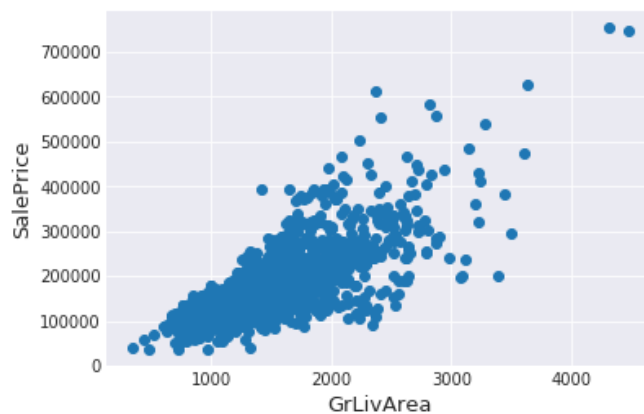
```
In [6]: fig, ax = plt.subplots()
ax.scatter(x = train['GrLivArea'], y = train['SalePrice'])
plt.ylabel('SalePrice', fontsize=13)
plt.xlabel('GrLivArea', fontsize=13)
plt.show()
```



We can see at the bottom right two with extremely large GrLivArea that are of a low price. These values are huge outliers. Therefore, we can safely delete them.

```
In [7]: #Deleting outliers
train = train.drop(train[(train['GrLivArea']>4000) & (train['SalePrice']<300000)].index)

#Check the graphic again
fig, ax = plt.subplots()
ax.scatter(train['GrLivArea'], train['SalePrice'])
plt.ylabel('SalePrice', fontsize=13)
plt.xlabel('GrLivArea', fontsize=13)
plt.show()
```



### Note :

Outliers removal is not always safe. We decided to delete these two as they are very huge and really bad ( extremely large areas for very low prices).

There are probably others outliers in the training data. However, removing all them may affect badly our models if ever there were also outliers in the test data. That's why , instead of removing them all, we will just manage to make some of our models robust on them. You can refer to the modelling part of this notebook for that.

## Target Variable

**SalePrice** is the variable we need to predict. So let's do some analysis on this variable first.

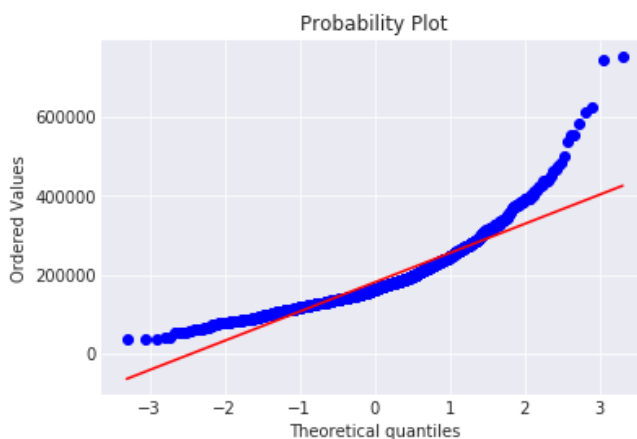
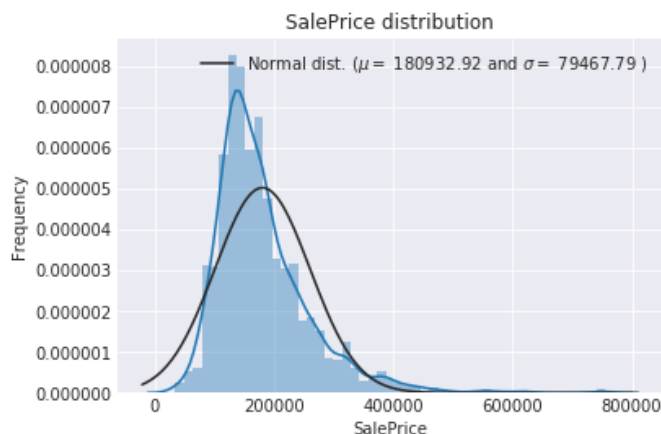
```
In [8]: sns.distplot(train['SalePrice'], fit=norm);

# 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')
plt.ylabel('Frequency')
plt.title('SalePrice distribution')

#Get also the QQ-plot
fig = plt.figure()
res = stats.probplot(train['SalePrice'], plot=plt)
plt.show()
```

mu = 180932.92 and sigma = 79467.79



The target variable is right skewed. As (linear) models love normally distributed data, we need to transform this variable and make it more normally distributed.

**Log-transformation of the target variable**

```
In [9]: #We use the numpy fuction log1p which applies log(1+x) to all elements of t
        #he column
        train["SalePrice"] = np.log1p(train["SalePrice"])

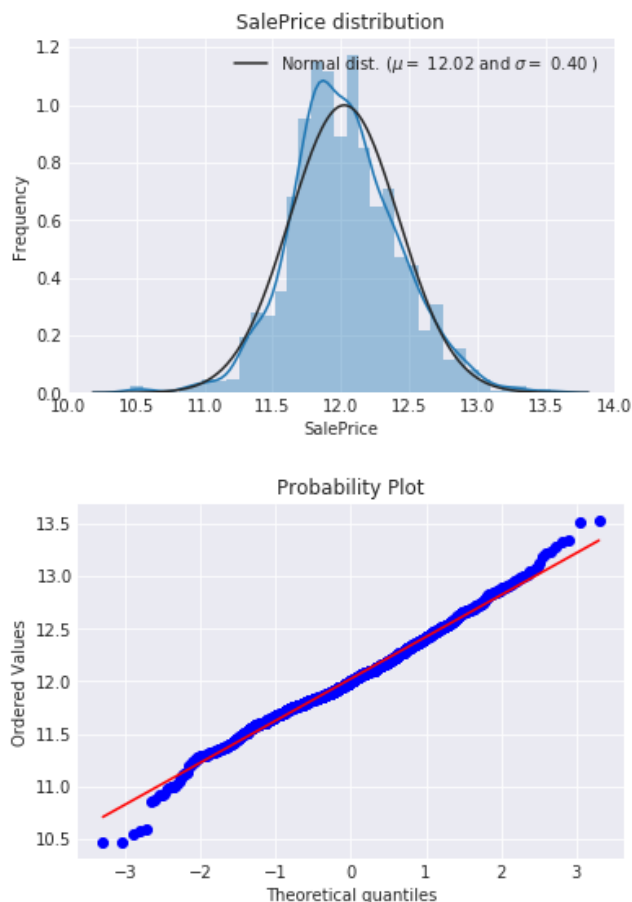
        #Check the new distribution
        sns.distplot(train['SalePrice'] , fit=norm);

        # 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')
        plt.ylabel('Frequency')
        plt.title('SalePrice distribution')

        #Get also the QQ-plot
        fig = plt.figure()
        res = stats.probplot(train['SalePrice'], plot=plt)
        plt.show()
```

mu = 12.02 and sigma = 0.40



The skew seems now corrected and the data appears more normally distributed.

## Features engineering

let's first concatenate the train and test data in the same dataframe

```
In [10]: ntrain = train.shape[0]
ntest = test.shape[0]
y_train = train.SalePrice.values
all_data = pd.concat((train, test)).reset_index(drop=True)
all_data.drop(['SalePrice'], axis=1, inplace=True)
print("all_data size is : {}".format(all_data.shape))

all_data size is : (2917, 79)
```

## Missing Data

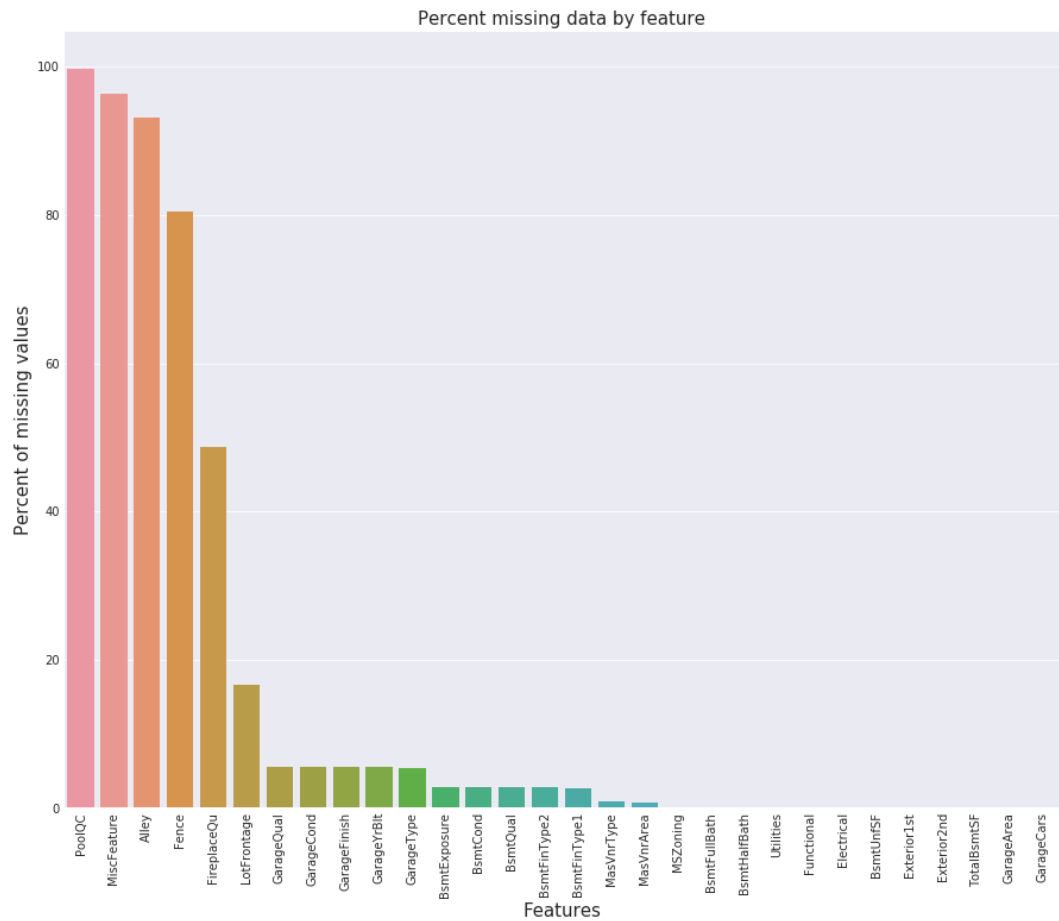
```
In [11]: all_data_na = (all_data.isnull().sum() / len(all_data)) * 100
all_data_na = all_data_na.drop(all_data_na[all_data_na == 0].index).sort_val
ues(ascending=False)[:30]
missing_data = pd.DataFrame({'Missing Ratio' :all_data_na})
missing_data.head(20)
```

Out[11]:

	Missing Ratio
PoolQC	99.691
MiscFeature	96.400
Alley	93.212
Fence	80.425
FireplaceQu	48.680
LotFrontage	16.661
GarageQual	5.451
GarageCond	5.451
GarageFinish	5.451
GarageYrBlt	5.451
GarageType	5.382
BsmtExposure	2.811
BsmtCond	2.811
BsmtQual	2.777
BsmtFinType2	2.743
BsmtFinType1	2.708
MasVnrType	0.823
MasVnrArea	0.788
MSZoning	0.137
BsmtFullBath	0.069

```
In [12]: f, ax = plt.subplots(figsize=(15, 12))
plt.xticks(rotation='90')
sns.barplot(x=all_data_na.index, y=all_data_na)
plt.xlabel('Features', fontsize=15)
plt.ylabel('Percent of missing values', fontsize=15)
plt.title('Percent missing data by feature', fontsize=15)
```

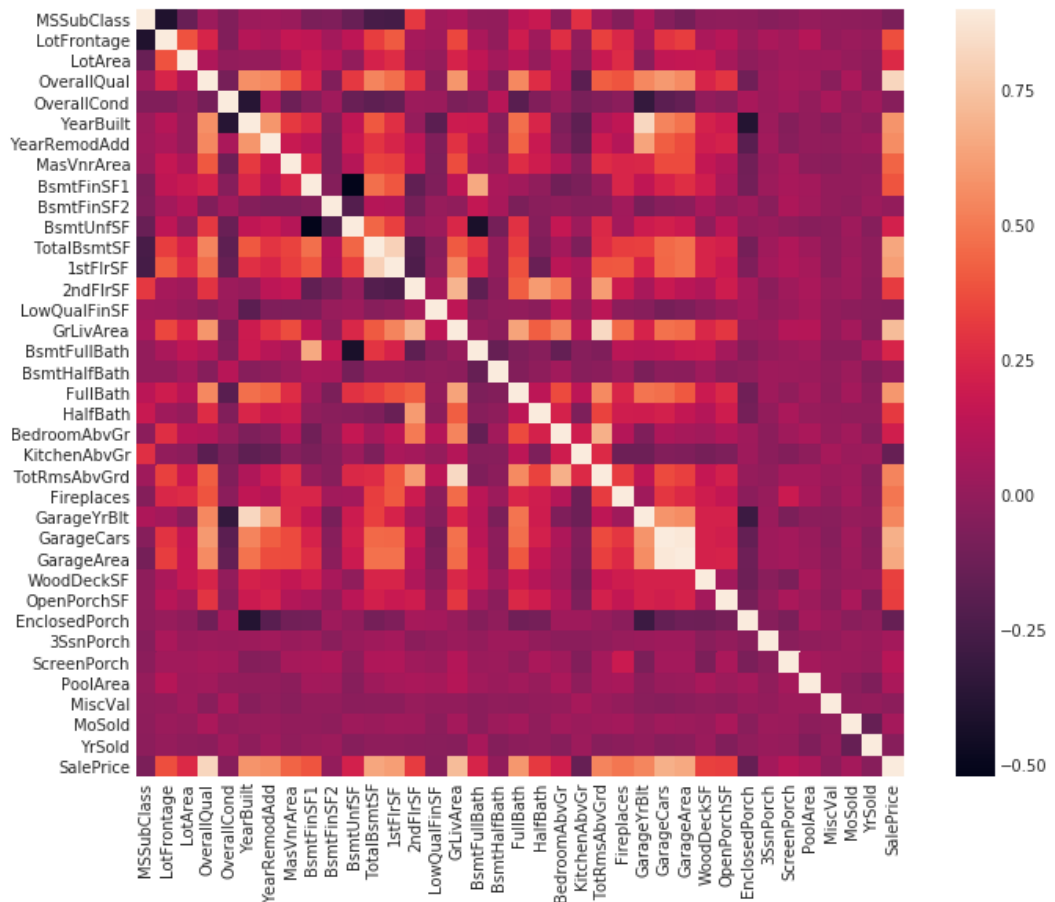
```
Out[12]: Text(0.5,1,'Percent missing data by feature')
```



## Data Correlation

```
In [13]: #Correlation map to see how features are correlated with SalePrice
corrmat = train.corr()
plt.subplots(figsize=(12,9))
sns.heatmap(corrmat, vmax=0.9, square=True)
```

```
Out[13]: <matplotlib.axes._subplots.AxesSubplot at 0x7efd7b454898>
```



## Imputing missing values

We impute them by proceeding sequentially through features with missing values

- **PoolQC** : data description says NA means "No Pool". That make sense, given the huge ratio of missing value (+99%) and majority of houses have no Pool at all in general.

```
In [14]: all_data["PoolQC"] = all_data["PoolQC"].fillna("None")
```

- **MiscFeature** : data description says NA means "no misc feature"

```
In [15]: all_data["MiscFeature"] = all_data["MiscFeature"].fillna("None")
```

- **Alley** : data description says NA means "no alley access"



```
In [16]: all_data["Alley"] = all_data["Alley"].fillna("None")
```

- **Fence** : data description says NA means "no fence"

```
In [17]: all_data["Fence"] = all_data["Fence"].fillna("None")
```

- **FireplaceQu** : data description says NA means "no fireplace"

```
In [18]: all_data["FireplaceQu"] = all_data["FireplaceQu"].fillna("None")
```

- **LotFrontage** : Since the area of each street connected to the house property most likely have a similar area to other houses in its neighborhood , we can **fill in missing values by the median LotFrontage of the neighborhood**.

```
In [19]: #Group by neighborhood and fill in missing value by the median LotFrontage of all the neighborhood
all_data["LotFrontage"] = all_data.groupby("Neighborhood")["LotFrontage"].transform(
    lambda x: x.fillna(x.median()))
```

- **GarageType, GarageFinish, GarageQual and GarageCond** : Replacing missing data with None

```
In [20]: for col in ('GarageType', 'GarageFinish', 'GarageQual', 'GarageCond'):
    all_data[col] = all_data[col].fillna('None')
```

- **GarageYrBlt, GarageArea and GarageCars** : Replacing missing data with 0 (Since No garage = no cars in such garage.)

```
In [21]: for col in ('GarageYrBlt', 'GarageArea', 'GarageCars'):
    all_data[col] = all_data[col].fillna(0)
```

- **BsmtFinSF1, BsmtFinSF2, BsmtUnfSF, TotalBsmtSF, BsmtFullBath and BsmtHalfBath** : missing values are likely zero for having no basement

```
In [22]: for col in ('BsmtFinSF1', 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', 'BsmtFullBath', 'BsmtHalfBath'):
    all_data[col] = all_data[col].fillna(0)
```

- **BsmtQual, BsmtCond, BsmtExposure, BsmtFinType1 and BsmtFinType2** : For all these categorical basement-related features, NaN means that there is no basement.

```
In [23]: for col in ('BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2'):
    all_data[col] = all_data[col].fillna('None')
```

- **MasVnrArea and MasVnrType** : NA most likely means no masonry veneer for these houses. We can fill 0 for the area and None for the type.

```
In [24]: all_data["MasVnrType"] = all_data["MasVnrType"].fillna("None")
all_data["MasVnrArea"] = all_data["MasVnrArea"].fillna(0)
```

- **MSZoning (The general zoning classification)** : 'RL' is by far the most common value. So we can fill in missing values with 'RL'

```
In [25]: all_data['MSZoning'] = all_data['MSZoning'].fillna(all_data['MSZoning'].mode()[0])
```

- **Utilities** : For this categorical feature all records are "AllPub", except for one "NoSeWa" and 2 NA . Since the house with 'NoSewa' is in the training set, **this feature won't help in predictive modelling**. We can then safely remove it.

```
In [26]: all_data = all_data.drop(['Utilities'], axis=1)
```

- **Functional** : data description says NA means typical

```
In [27]: all_data["Functional"] = all_data["Functional"].fillna("Typ")
```

- **Electrical** : It has one NA value. Since this feature has mostly 'SBrkr', we can set that for the missing value.

```
In [28]: all_data['Electrical'] = all_data['Electrical'].fillna(all_data['Electrical'].mode()[0])
```

- **KitchenQual**: Only one NA value, and same as Electrical, we set 'TA' (which is the most frequent) for the missing value in KitchenQual.

```
In [29]: all_data['KitchenQual'] = all_data['KitchenQual'].fillna(all_data['KitchenQual'].mode()[0])
```

- **Exterior1st and Exterior2nd** : Again Both Exterior 1 & 2 have only one missing value. We will just substitute in the most common string

```
In [30]: all_data['Exterior1st'] = all_data['Exterior1st'].fillna(all_data['Exterior1st'].mode()[0])
all_data['Exterior2nd'] = all_data['Exterior2nd'].fillna(all_data['Exterior2nd'].mode()[0])
```

- **SaleType** : Fill in again with most frequent which is "WD"

```
In [31]: all_data['SaleType'] = all_data['SaleType'].fillna(all_data['SaleType'].mode()[0])
```

- **MSSubClass** : Na most likely means No building class. We can replace missing values with None

```
In [32]: all_data['MSSubClass'] = all_data['MSSubClass'].fillna("None")
```

Is there any remaining missing value ?

```
In [33]: #Check remaining missing values if any
all_data_na = (all_data.isnull().sum() / len(all_data)) * 100
all_data_na = all_data_na.drop(all_data_na[all_data_na == 0].index).sort_val
ues(ascending=False)
missing_data = pd.DataFrame({'Missing Ratio' :all_data_na})
missing_data.head()
```

```
Out[33]:
```

Missing Ratio
---------------

It remains no missing value.

## More features engineering

Transforming some numerical variables that are really categorical

```
In [34]: #MSSubClass=The building class
all_data['MSSubClass'] = all_data['MSSubClass'].apply(str)

#Changing OverallCond into a categorical variable
all_data['OverallCond'] = all_data['OverallCond'].astype(str)

#Year and month sold are transformed into categorical features.
all_data['YrSold'] = all_data['YrSold'].astype(str)
all_data['MoSold'] = all_data['MoSold'].astype(str)
```

Label Encoding some categorical variables that may contain information in their ordering set

```
In [35]: from sklearn.preprocessing import LabelEncoder
cols = ('FireplaceQu', 'BsmtQual', 'BsmtCond', 'GarageQual', 'GarageCond',
        'ExterQual', 'ExterCond', 'HeatingQC', 'PoolQC', 'KitchenQual', 'Bsmt
        FinType1',
        'BsmtFinType2', 'Functional', 'Fence', 'BsmtExposure', 'GarageFinish
        ', 'LandSlope',
        'LotShape', 'PavedDrive', 'Street', 'Alley', 'CentralAir', 'MSSubCla
        ss', 'OverallCond',
        'YrSold', 'MoSold')
# process columns, apply LabelEncoder to categorical features
for c in cols:
    lbl = LabelEncoder()
    lbl.fit(list(all_data[c].values))
    all_data[c] = lbl.transform(list(all_data[c].values))

# shape
print('Shape all_data: {}'.format(all_data.shape))
```

```
Shape all_data: (2917, 78)
```

Adding one more important feature

Since area related features are very important to determine house prices, we add one more feature which is the total area of basement, first and second floor areas of each house

```
In [36]: # Adding total sqfootage feature
all_data['TotalSF'] = all_data['TotalBsmtSF'] + all_data['1stFlrSF'] + all_data['2ndFlrSF']
```

### Skewed features

```
In [37]: numeric_feats = all_data.dtypes[all_data.dtypes != "object"].index

# Check the skew of all numerical features
skewed_feats = all_data[numeric_feats].apply(lambda x: skew(x.dropna())).sort_values(ascending=False)
print("\nSkew in numerical features: \n")
skewness = pd.DataFrame({'Skew' :skewed_feats})
skewness.head(10)
```

Skew in numerical features:

Out[37]:

	Skew
MiscVal	21.940
PoolArea	17.689
LotArea	13.109
LowQualFinSF	12.085
3SsnPorch	11.372
LandSlope	4.973
KitchenAbvGr	4.301
BsmtFinSF2	4.145
EnclosedPorch	4.002
ScreenPorch	3.945

### Box Cox Transformation of (highly) skewed features

We use the scipy function `boxcox1p` which computes the Box-Cox transformation of  $1 + x$ .

Note that setting  $\lambda = 0$  is equivalent to `log1p` used above for the target variable.

See [this page \(http://onlinestatbook.com/2/transformations/box-cox.html\)](http://onlinestatbook.com/2/transformations/box-cox.html) for more details on Box Cox Transformation as well as [the scipy function's page \(https://docs.scipy.org/doc/scipy-0.19.0/reference/generated/scipy.special.boxcox1p.html\)](https://docs.scipy.org/doc/scipy-0.19.0/reference/generated/scipy.special.boxcox1p.html)

```
In [38]: skewness = skewness[abs(skewness) > 0.75]
print("There are {} skewed numerical features to Box Cox transform".format(
    skewness.shape[0]))

from scipy.special import boxcox1p
skewed_features = skewness.index
lam = 0.15
for feat in skewed_features:
    #all_data[feat] += 1
    all_data[feat] = boxcox1p(all_data[feat], lam)

#all_data[skewed_features] = np.log1p(all_data[skewed_features])
```

There are 59 skewed numerical features to Box Cox transform

### Getting dummy categorical features

```
In [39]: all_data = pd.get_dummies(all_data)
print(all_data.shape)

(2917, 220)
```

Getting the new train and test sets.

```
In [40]: train = all_data[:ntrain]
test = all_data[ntrain:]
```

## Modelling

### Import libraries

```
In [41]: from sklearn.linear_model import ElasticNet, Lasso, BayesianRidge, LassoLar
sIC
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.kernel_ridge import KernelRidge
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import RobustScaler
from sklearn.base import BaseEstimator, TransformerMixin, RegressorMixin, clone
one
from sklearn.model_selection import KFold, cross_val_score, train_test_split
from sklearn.metrics import mean_squared_error
import xgboost as xgb
import lightgbm as lgb
```

### Define a cross validation strategy

We use the **cross\_val\_score** function of Sklearn. However this function has not a shuffle attribut, we add then one line of code, in order to shuffle the dataset prior to cross-validation

```
In [42]: #Validation function
n_folds = 5

def rmsle_cv(model):
    kf = KFold(n_folds, shuffle=True, random_state=42).get_n_splits(train.values)
    rmse= np.sqrt(-cross_val_score(model, train.values, y_train, scoring="negative_mean_squared_error", cv = kf))
    return(rmse)
```

## Base models

- **LASSO Regression :**

This model may be very sensitive to outliers. So we need to made it more robust on them. For that we use the sklearn's **RobustScaler()** method on pipeline

```
In [43]: lasso = make_pipeline(RobustScaler(), Lasso(alpha =0.0005, random_state=1))
```

- **Elastic Net Regression :**

again made robust to outliers

```
In [44]: ENet = make_pipeline(RobustScaler(), ElasticNet(alpha=0.0005, l1_ratio=.9, random_state=3))
```

- **Kernel Ridge Regression :**

```
In [45]: KRR = KernelRidge(alpha=0.6, kernel='polynomial', degree=2, coef0=2.5)
```

- **Gradient Boosting Regression :**

With **huber** loss that makes it robust to outliers

```
In [46]: GBoost = GradientBoostingRegressor(n_estimators=3000, learning_rate=0.05,
                                             max_depth=4, max_features='sqrt',
                                             min_samples_leaf=15, min_samples_split=10,
                                             loss='huber', random_state =5)
```

- **XGBoost :**

```
In [47]: model_xgb = xgb.XGBRegressor(colsample_bytree=0.4603, gamma=0.0468,
                                       learning_rate=0.05, max_depth=3,
                                       min_child_weight=1.7817, n_estimators=2200,
                                       reg_alpha=0.4640, reg_lambda=0.8571,
                                       subsample=0.5213, silent=1,
                                       random_state =7, nthread = -1)
```

- **LightGBM :**

```
In [48]: model_lgb = lgb.LGBMRegressor(objective='regression', num_leaves=5,
                                         learning_rate=0.05, n_estimators=720,
                                         max_bin = 55, bagging_fraction = 0.8,
                                         bagging_freq = 5, feature_fraction = 0.2319,
                                         feature_fraction_seed=9, bagging_seed=9,
                                         min_data_in_leaf =6, min_sum_hessian_in_leaf =
11)
```

## Base models scores

Let's see how these base models perform on the data by evaluating the cross-validation rmsle error

```
In [49]: score = rmsle_cv(lasso)
print("\nLasso score: {:.4f} ({:.4f})\n".format(score.mean(), score.std()))

Lasso score: 0.1115 (0.0074)
```

```
In [50]: score = rmsle_cv(ENet)
print("ElasticNet score: {:.4f} ({:.4f})\n".format(score.mean(), score.std
()))

ElasticNet score: 0.1116 (0.0074)
```

```
In [51]: score = rmsle_cv(KRR)
print("Kernel Ridge score: {:.4f} ({:.4f})\n".format(score.mean(), score.std
()))

Kernel Ridge score: 0.1153 (0.0075)
```

```
In [52]: score = rmsle_cv(GBoost)
print("Gradient Boosting score: {:.4f} ({:.4f})\n".format(score.mean(), scor
e.std()))

Gradient Boosting score: 0.1177 (0.0080)
```

```
In [53]: score = rmsle_cv(model_xgb)
print("Xgboost score: {:.4f} ({:.4f})\n".format(score.mean(), score.std()))

Xgboost score: 0.1161 (0.0079)
```

```
In [54]: score = rmsle_cv(model_lgb)
print("LGBM score: {:.4f} ({:.4f})\n" .format(score.mean(), score.std()))

LGBM score: 0.1157 (0.0067)
```

## Stacking models

### Simplest Stacking approach : Averaging base models

We begin with this simple approach of averaging base models. We build a new **class** to extend scikit-learn with our model and also to leverage encapsulation and code reuse ([inheritance \(https://en.wikipedia.org/wiki/Inheritance\\_\(object-oriented\\_programming\)\)](https://en.wikipedia.org/wiki/Inheritance_(object-oriented_programming))))

### Averaged base models class

```
In [55]: class AveragingModels(BaseEstimator, RegressorMixin, TransformerMixin):
def __init__(self, models):
    self.models = models

    # we define clones of the original models to fit the data in
def fit(self, X, y):
    self.models_ = [clone(x) for x in self.models]

    # Train cloned base models
    for model in self.models_:
        model.fit(X, y)

    return self

    #Now we do the predictions for cloned models and average them
def predict(self, X):
    predictions = np.column_stack([
        model.predict(X) for model in self.models_
    ])
    return np.mean(predictions, axis=1)
```

### Averaged base models score

We just average four models here **ENet**, **GBoost**, **KRR** and **lasso**. Of course we could easily add more models in the mix.

```
In [56]: averaged_models = AveragingModels(models = (ENet, GBoost, KRR, lasso))

score = rmsle_cv(averaged_models)
print(" Averaged base models score: {:.4f} ({:.4f})\n".format(score.mean(),
score.std()))
```

Averaged base models score: 0.1091 (0.0075)

Wow ! It seems even the simplest stacking approach really improve the score . This encourages us to go further and explore a less simple stacking approach.

### Less simple Stacking : Adding a Meta-model



In this approach, we add a meta-model on averaged base models and use the out-of-folds predictions of these base models to train our meta-model.

The procedure, for the training part, may be described as follows:

1. Split the total training set into two disjoint sets (here **train** and **holdout** )
2. Train several base models on the first part (**train**)
3. Test these base models on the second part (**holdout**)
4. Use the predictions from 3) (called out-of-folds predictions) as the inputs, and the correct responses (target variable) as the outputs to train a higher level learner called **meta-model**.

The first three steps are done iteratively . If we take for example a 5-fold stacking , we first split the training data into 5 folds. Then we will do 5 iterations. In each iteration, we train every base model on 4 folds and predict on the remaining fold (holdout fold).

So, we will be sure, after 5 iterations , that the entire data is used to get out-of-folds predictions that we will then use as new feature to train our meta-model in the step 4.

For the prediction part , We average the predictions of all base models on the test data and used them as **meta-features** on which, the final prediction is done with the meta-model.

Faron

(Image taken from [Faron \(https://www.kaggle.com/getting-started/18153#post103381\)](https://www.kaggle.com/getting-started/18153#post103381))

kaz

Gif taken from [KazAnova's interview \(http://blog.kaggle.com/2017/06/15/stacking-made-easy-an-introduction-to-stacknet-by-competitions-grandmaster-marios-michailidis-kazanova/\)](http://blog.kaggle.com/2017/06/15/stacking-made-easy-an-introduction-to-stacknet-by-competitions-grandmaster-marios-michailidis-kazanova/)

On this gif, the base models are algorithms 0, 1, 2 and the meta-model is algorithm 3. The entire training dataset is A+B (target variable y known) that we can split into train part (A) and holdout part (B). And the test dataset is C.

B1 (which is the prediction from the holdout part) is the new feature used to train the meta-model 3 and C1 (which is the prediction from the test dataset) is the meta-feature on which the final prediction is done.

### Stacking averaged Models Class

```
In [57]: class StackingAveragedModels(BaseEstimator, RegressorMixin, TransformerMixin):
def __init__(self, base_models, meta_model, n_folds=5):
    self.base_models = base_models
    self.meta_model = meta_model
    self.n_folds = n_folds

    # We again fit the data on clones of the original models
def fit(self, X, y):
    self.base_models_ = [list() for x in self.base_models]
    self.meta_model_ = clone(self.meta_model)
    kfold = KFold(n_splits=self.n_folds, shuffle=True, random_state=156)

    # Train cloned base models then create out-of-fold predictions
    # that are needed to train the cloned meta-model
    out_of_fold_predictions = np.zeros((X.shape[0], len(self.base_model
s)))

    for i, model in enumerate(self.base_models):
        for train_index, holdout_index in kfold.split(X, y):
            instance = clone(model)
            self.base_models_[i].append(instance)
            instance.fit(X[train_index], y[train_index])
            y_pred = instance.predict(X[holdout_index])
            out_of_fold_predictions[holdout_index, i] = y_pred

    # Now train the cloned meta-model using the out-of-fold predictions
    # as new feature
    self.meta_model_.fit(out_of_fold_predictions, y)
    return self

    # Do the predictions of all base models on the test data and use the aver
    # aged predictions as
    # meta-features for the final prediction which is done by the meta-model
def predict(self, X):
    meta_features = np.column_stack([
        np.column_stack([model.predict(X) for model in base_models]).mean
n(axis=1)
        for base_models in self.base_models_ ])
    return self.meta_model_.predict(meta_features)
```

### Stacking Averaged models Score

To make the two approaches comparable (by using the same number of models) , we just average **Enet KRR** and **Gboost**, then we add **lasso as meta-model**.

```
In [58]: stacked_averaged_models = StackingAveragedModels(base_models = (ENet, GBoos
t, KRR),
                                                         meta_model = lasso)

score = rmsle_cv(stacked_averaged_models)
print("Stacking Averaged models score: {:.4f} ({:.4f})".format(score.mean(),
score.std()))

Stacking Averaged models score: 0.1085 (0.0074)
```

We get again a better score by adding a meta learner

## Ensembling StackedRegressor, XGBoost and LightGBM

We add **XGBoost** and **LightGBM** to the **StackedRegressor** defined previously.

We first define a rmsle evaluation function

```
In [59]: def rmsle(y, y_pred):  
         return np.sqrt(mean_squared_error(y, y_pred))
```

## Final Training and Prediction

**StackedRegressor:**

```
In [60]: stacked_averaged_models.fit(train.values, y_train)  
stacked_train_pred = stacked_averaged_models.predict(train.values)  
stacked_pred = np.expml(stacked_averaged_models.predict(test.values))  
print(rmsle(y_train, stacked_train_pred))  
  
0.0781571937916
```

**XGBoost:**

```
In [61]: model_xgb.fit(train, y_train)  
xgb_train_pred = model_xgb.predict(train)  
xgb_pred = np.expml(model_xgb.predict(test))  
print(rmsle(y_train, xgb_train_pred))  
  
0.0785165142425
```

**LightGBM:**

```
In [62]: model_lgb.fit(train, y_train)  
lgb_train_pred = model_lgb.predict(train)  
lgb_pred = np.expml(model_lgb.predict(test.values))  
print(rmsle(y_train, lgb_train_pred))  
  
0.0716757468834
```

```
In [63]: '''RMSE on the entire Train data when averaging'''  
  
print('RMSLE score on train data:')  
print(rmsle(y_train, stacked_train_pred*0.70 +  
           xgb_train_pred*0.15 + lgb_train_pred*0.15 ))  
  
RMSLE score on train data:  
0.0752190464543
```

**Ensemble prediction:**

```
In [64]: ensemble = stacked_pred*0.70 + xgb_pred*0.15 + lgb_pred*0.15
```

**Submission**

```
In [65]: sub = pd.DataFrame()  
sub['Id'] = test_ID  
sub['SalePrice'] = ensemble  
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 :-)

```

# This script creates a ton of features, then trains an XGBoost regressor
# and a Lasso regressor, and combines their predictions.
#
# It borrows ideas from lots of other people's scripts, including:
# https://www.kaggle.com/klyusba/house-prices-advanced-regression-techniques/lasso-
# https://www.kaggle.com/juliencs/house-prices-advanced-regression-techniques/a-
# https://www.kaggle.com/apapiu/house-prices-advanced-regression-techniques/
# regularized-linear-models
# but I probably forgot to mention a few. ;-)

import numpy as np
import pandas as pd

# The error metric: RMSE on the log of the sale prices.
from sklearn.metrics import mean_squared_error
def rmse(y_true, y_pred):
    return np.sqrt(mean_squared_error(y_true, y_pred))

# Load the data.
train_df = pd.read_csv("../input/train.csv")
test_df = pd.read_csv("../input/test.csv")

# There are a few houses with more than 4000 sq ft living area that are
# outliers, so we drop them from the training data. (There is also one in
# the test set but we obviously can't drop that one.)
train_df.drop(train_df[train_df["GrLivArea"] > 4000].index, inplace=True)

# The test example with ID 666 has GarageArea, GarageCars, and GarageType
# but none of the other fields, so use the mode and median to fill them in.
test_df.loc[666, "GarageQual"] = "TA"
test_df.loc[666, "GarageCond"] = "TA"
test_df.loc[666, "GarageFinish"] = "Unf"
test_df.loc[666, "GarageYrBlt"] = "1980"

# The test example 1116 only has GarageType but no other information. We'll
# assume it does not have a garage.
test_df.loc[1116, "GarageType"] = np.nan

# For imputing missing values: fill in missing LotFrontage values by the median
# LotFrontage of the neighborhood.
lot_frontage_by_neighborhood =
train_df["LotFrontage"].groupby(train_df["Neighborhood"])

# Used to convert categorical features into ordinal numbers.
# (There's probably an easier way to do this, but it works.)
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()

def factorize(df, factor_df, column, fill_na=None):
    factor_df[column] = df[column]
    if fill_na is not None:
        factor_df[column].fillna(fill_na, inplace=True)
    le.fit(factor_df[column].unique())
    factor_df[column] = le.transform(factor_df[column])
    return factor_df

# Combine all the (numerical) features into one big DataFrame. We don't add
# the one-hot encoded variables here yet, that happens later on.
def munge(df):
    all_df = pd.DataFrame(index = df.index)

    all_df["LotFrontage"] = df["LotFrontage"]
    for key, group in lot_frontage_by_neighborhood:
        idx = (df["Neighborhood"] == key) & (df["LotFrontage"].isnull())

```

```

    all_df.loc[idx, "LotFrontage"] = group.median()

all_df["LotArea"] = df["LotArea"]

all_df["MasVnrArea"] = df["MasVnrArea"]
all_df["MasVnrArea"].fillna(0, inplace=True)

all_df["BsmtFinSF1"] = df["BsmtFinSF1"]
all_df["BsmtFinSF1"].fillna(0, inplace=True)

all_df["BsmtFinSF2"] = df["BsmtFinSF2"]
all_df["BsmtFinSF2"].fillna(0, inplace=True)

all_df["BsmtUnfSF"] = df["BsmtUnfSF"]
all_df["BsmtUnfSF"].fillna(0, inplace=True)

all_df["TotalBsmtSF"] = df["TotalBsmtSF"]
all_df["TotalBsmtSF"].fillna(0, inplace=True)

all_df["1stFlrSF"] = df["1stFlrSF"]
all_df["2ndFlrSF"] = df["2ndFlrSF"]
all_df["GrLivArea"] = df["GrLivArea"]

all_df["GarageArea"] = df["GarageArea"]
all_df["GarageArea"].fillna(0, inplace=True)

all_df["WoodDeckSF"] = df["WoodDeckSF"]
all_df["OpenPorchSF"] = df["OpenPorchSF"]
all_df["EnclosedPorch"] = df["EnclosedPorch"]
all_df["3SsnPorch"] = df["3SsnPorch"]
all_df["ScreenPorch"] = df["ScreenPorch"]

all_df["BsmtFullBath"] = df["BsmtFullBath"]
all_df["BsmtFullBath"].fillna(0, inplace=True)

all_df["BsmtHalfBath"] = df["BsmtHalfBath"]
all_df["BsmtHalfBath"].fillna(0, inplace=True)

all_df["FullBath"] = df["FullBath"]
all_df["HalfBath"] = df["HalfBath"]
all_df["BedroomAbvGr"] = df["BedroomAbvGr"]
all_df["KitchenAbvGr"] = df["KitchenAbvGr"]
all_df["TotRmsAbvGrd"] = df["TotRmsAbvGrd"]
all_df["Fireplaces"] = df["Fireplaces"]

all_df["GarageCars"] = df["GarageCars"]
all_df["GarageCars"].fillna(0, inplace=True)

all_df["CentralAir"] = (df["CentralAir"] == "Y") * 1.0

all_df["OverallQual"] = df["OverallQual"]
all_df["OverallCond"] = df["OverallCond"]

# Quality measurements are stored as text but we can convert them to
# numbers where a higher number means higher quality.

qual_dict = {None: 0, "Po": 1, "Fa": 2, "TA": 3, "Gd": 4, "Ex": 5}
all_df["ExterQual"] = df["ExterQual"].map(qual_dict).astype(int)
all_df["ExterCond"] = df["ExterCond"].map(qual_dict).astype(int)
all_df["BsmtQual"] = df["BsmtQual"].map(qual_dict).astype(int)
all_df["BsmtCond"] = df["BsmtCond"].map(qual_dict).astype(int)
all_df["HeatingQC"] = df["HeatingQC"].map(qual_dict).astype(int)
all_df["KitchenQual"] = df["KitchenQual"].map(qual_dict).astype(int)
all_df["FireplaceQu"] = df["FireplaceQu"].map(qual_dict).astype(int)
all_df["GarageQual"] = df["GarageQual"].map(qual_dict).astype(int)
all_df["GarageCond"] = df["GarageCond"].map(qual_dict).astype(int)

```

```

all_df["BsmtExposure"] = df["BsmtExposure"].map(
    {None: 0, "No": 1, "Mn": 2, "Av": 3, "Gd": 4}).astype(int)

bsmt_fin_dict = {None: 0, "Unf": 1, "LwQ": 2, "Rec": 3, "BLQ": 4, "ALQ": 5,
"GLQ": 6}
all_df["BsmtFinType1"] = df["BsmtFinType1"].map(bsmt_fin_dict).astype(int)
all_df["BsmtFinType2"] = df["BsmtFinType2"].map(bsmt_fin_dict).astype(int)

all_df["Functional"] = df["Functional"].map(
    {None: 0, "Sal": 1, "Sev": 2, "Maj2": 3, "Maj1": 4,
    "Mod": 5, "Min2": 6, "Min1": 7, "Typ": 8}).astype(int)

all_df["GarageFinish"] = df["GarageFinish"].map(
    {None: 0, "Unf": 1, "RFn": 2, "Fin": 3}).astype(int)

all_df["Fence"] = df["Fence"].map(
    {None: 0, "MnWw": 1, "GdWo": 2, "MnPrv": 3, "GdPrv": 4}).astype(int)

all_df["YearBuilt"] = df["YearBuilt"]
all_df["YearRemodAdd"] = df["YearRemodAdd"]

all_df["GarageYrBlt"] = df["GarageYrBlt"]
all_df["GarageYrBlt"].fillna(0.0, inplace=True)

all_df["MoSold"] = df["MoSold"]
all_df["YrSold"] = df["YrSold"]

all_df["LowQualFinSF"] = df["LowQualFinSF"]
all_df["MiscVal"] = df["MiscVal"]

all_df["PoolQC"] = df["PoolQC"].map(qual_dict).astype(int)

all_df["PoolArea"] = df["PoolArea"]
all_df["PoolArea"].fillna(0, inplace=True)

# Add categorical features as numbers too. It seems to help a bit.
all_df = factorize(df, all_df, "MSSubClass")
all_df = factorize(df, all_df, "MSZoning", "RL")
all_df = factorize(df, all_df, "LotConfig")
all_df = factorize(df, all_df, "Neighborhood")
all_df = factorize(df, all_df, "Condition1")
all_df = factorize(df, all_df, "BldgType")
all_df = factorize(df, all_df, "HouseStyle")
all_df = factorize(df, all_df, "RoofStyle")
all_df = factorize(df, all_df, "Exterior1st", "Other")
all_df = factorize(df, all_df, "Exterior2nd", "Other")
all_df = factorize(df, all_df, "MasVnrType", "None")
all_df = factorize(df, all_df, "Foundation")
all_df = factorize(df, all_df, "SaleType", "Oth")
all_df = factorize(df, all_df, "SaleCondition")

# IR2 and IR3 don't appear that often, so just make a distinction
# between regular and irregular.
all_df["IsRegularLotShape"] = (df["LotShape"] == "Reg") * 1

# Most properties are level; bin the other possibilities together
# as "not level".
all_df["IsLandLevel"] = (df["LandContour"] == "Lvl") * 1

# Most land slopes are gentle; treat the others as "not gentle".
all_df["IsLandSlopeGentle"] = (df["LandSlope"] == "Gtl") * 1

# Most properties use standard circuit breakers.
all_df["IsElectricalSBrkr"] = (df["Electrical"] == "SBrkr") * 1

```

```

# About 2/3rd have an attached garage.
all_df["IsGarageDetached"] = (df["GarageType"] == "Detchd") * 1

# Most have a paved drive. Treat dirt/gravel and partial pavement
# as "not paved".
all_df["IsPavedDrive"] = (df["PavedDrive"] == "Y") * 1

# The only interesting "misc. feature" is the presence of a shed.
all_df["HasShed"] = (df["MiscFeature"] == "Shed") * 1.

# If YearRemodAdd != YearBuilt, then a remodeling took place at some point.
all_df["Remodeled"] = (all_df["YearRemodAdd"] != all_df["YearBuilt"]) * 1

# Did a remodeling happen in the year the house was sold?
all_df["RecentRemodel"] = (all_df["YearRemodAdd"] == all_df["YrSold"]) * 1

# Was this house sold in the year it was built?
all_df["VeryNewHouse"] = (all_df["YearBuilt"] == all_df["YrSold"]) * 1

all_df["Has2ndFloor"] = (all_df["2ndFlrSF"] == 0) * 1
all_df["HasMasVnr"] = (all_df["MasVnrArea"] == 0) * 1
all_df["HasWoodDeck"] = (all_df["WoodDeckSF"] == 0) * 1
all_df["HasOpenPorch"] = (all_df["OpenPorchSF"] == 0) * 1
all_df["HasEnclosedPorch"] = (all_df["EnclosedPorch"] == 0) * 1
all_df["Has3SsnPorch"] = (all_df["3SsnPorch"] == 0) * 1
all_df["HasScreenPorch"] = (all_df["ScreenPorch"] == 0) * 1

# These features actually lower the score a little.
# all_df["HasBasement"] = df["BsmtQual"].isnull() * 1
# all_df["HasGarage"] = df["GarageQual"].isnull() * 1
# all_df["HasFireplace"] = df["FireplaceQu"].isnull() * 1
# all_df["HasFence"] = df["Fence"].isnull() * 1

# Months with the largest number of deals may be significant.
all_df["HighSeason"] = df["MoSold"].replace(
    {1: 0, 2: 0, 3: 0, 4: 1, 5: 1, 6: 1, 7: 1, 8: 0, 9: 0, 10: 0, 11: 0, 12:
0})

all_df["NewerDwelling"] = df["MSSubClass"].replace(
    {20: 1, 30: 0, 40: 0, 45: 0, 50: 0, 60: 1, 70: 0, 75: 0, 80: 0, 85: 0,
    90: 0, 120: 1, 150: 0, 160: 0, 180: 0, 190: 0})

all_df.loc[df.Neighborhood == 'NridgHt', "Neighborhood_Good"] = 1
all_df.loc[df.Neighborhood == 'Crawfor', "Neighborhood_Good"] = 1
all_df.loc[df.Neighborhood == 'StoneBr', "Neighborhood_Good"] = 1
all_df.loc[df.Neighborhood == 'Somerst', "Neighborhood_Good"] = 1
all_df.loc[df.Neighborhood == 'NoRidge', "Neighborhood_Good"] = 1
all_df["Neighborhood_Good"].fillna(0, inplace=True)

all_df["SaleCondition_PriceDown"] = df.SaleCondition.replace(
    {'Abnorml': 1, 'Alloca': 1, 'AdjLand': 1, 'Family': 1, 'Normal': 0,
'Partial': 0})

# House completed before sale or not
all_df["BoughtOffPlan"] = df.SaleCondition.replace(
    {"Abnorml" : 0, "Alloca" : 0, "AdjLand" : 0, "Family" : 0, "Normal" : 0,
"Partial" : 1})

all_df["BadHeating"] = df.HeatingQC.replace(
    {'Ex': 0, 'Gd': 0, 'TA': 0, 'Fa': 1, 'Po': 1})

area_cols = ['LotFrontage', 'LotArea', 'MasVnrArea', 'BsmtFinSF1',
'BsmtFinSF2', 'BsmtUnfSF',
'TotalBsmtSF', '1stFlrSF', '2ndFlrSF', 'GrLivArea', 'GarageArea',
'WoodDeckSF',
'OpenPorchSF', 'EnclosedPorch', '3SsnPorch', 'ScreenPorch',

```



```

'LowQualFinSF', 'PoolArea' ]
all_df["TotalArea"] = all_df[area_cols].sum(axis=1)

all_df["TotalArea1st2nd"] = all_df["1stFlrSF"] + all_df["2ndFlrSF"]

all_df["Age"] = 2010 - all_df["YearBuilt"]
all_df["TimeSinceSold"] = 2010 - all_df["YrSold"]

all_df["SeasonSold"] = all_df["MoSold"].map({12:0, 1:0, 2:0, 3:1, 4:1, 5:1,
                                             6:2, 7:2, 8:2, 9:3, 10:3,
11:3}).astype(int)

all_df["YearsSinceRemodel"] = all_df["YrSold"] - all_df["YearRemodAdd"]

# Simplifications of existing features into bad/average/good.
all_df["SimplOverallQual"] = all_df.OverallQual.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2, 6 : 2, 7 : 3, 8 : 3, 9 : 3, 10 : 3})
all_df["SimplOverallCond"] = all_df.OverallCond.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2, 6 : 2, 7 : 3, 8 : 3, 9 : 3, 10 : 3})
all_df["SimplPoolQC"] = all_df.PoolQC.replace(
    {1 : 1, 2 : 1, 3 : 2, 4 : 2})
all_df["SimplGarageCond"] = all_df.GarageCond.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplGarageQual"] = all_df.GarageQual.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplFireplaceQu"] = all_df.FireplaceQu.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplFireplaceQu"] = all_df.FireplaceQu.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplFunctional"] = all_df.Functional.replace(
    {1 : 1, 2 : 1, 3 : 2, 4 : 2, 5 : 3, 6 : 3, 7 : 3, 8 : 4})
all_df["SimplKitchenQual"] = all_df.KitchenQual.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplHeatingQC"] = all_df.HeatingQC.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplBsmtFinType1"] = all_df.BsmtFinType1.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2, 6 : 2})
all_df["SimplBsmtFinType2"] = all_df.BsmtFinType2.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2, 6 : 2})
all_df["SimplBsmtCond"] = all_df.BsmtCond.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplBsmtQual"] = all_df.BsmtQual.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplExterCond"] = all_df.ExterCond.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})
all_df["SimplExterQual"] = all_df.ExterQual.replace(
    {1 : 1, 2 : 1, 3 : 1, 4 : 2, 5 : 2})

# Bin by neighborhood (a little arbitrarily). Values were computed by:
#
train_df["SalePrice"].groupby(train_df["Neighborhood"]).median().sort_values()
neighborhood_map = {
    "MeadowV" : 0, # 88000
    "IDOTRR" : 1, # 103000
    "BrDale" : 1, # 106000
    "OldTown" : 1, # 119000
    "Edwards" : 1, # 119500
    "BrkSide" : 1, # 124300
    "Sawyer" : 1, # 135000
    "Blueste" : 1, # 137500
    "SWISU" : 2, # 139500
    "NAmes" : 2, # 140000
    "NPkVill" : 2, # 146000
    "Mitchel" : 2, # 153500
    "SawyerW" : 2, # 179900
    "Gilbert" : 2, # 181000

```

```

        "NWAmes" : 2, # 182900
        "Blmngtn" : 2, # 191000
        "CollgCr" : 2, # 197200
        "ClearCr" : 3, # 200250
        "Crawfor" : 3, # 200624
        "Veenker" : 3, # 218000
        "Somerst" : 3, # 225500
        "Timber" : 3, # 228475
        "StoneBr" : 4, # 278000
        "NoRidge" : 4, # 290000
        "NridgHt" : 4, # 315000
    }

    all_df["NeighborhoodBin"] = df["Neighborhood"].map(neighborhood_map)
    return all_df

train_df_munged = munge(train_df)
test_df_munged = munge(test_df)

print(train_df_munged.shape)
print(test_df_munged.shape)

# Copy NeighborhoodBin into a temporary DataFrame because we want to use the
# unscaled version later on (to one-hot encode it).
neighborhood_bin_train = pd.DataFrame(index = train_df.index)
neighborhood_bin_train["NeighborhoodBin"] = train_df_munged["NeighborhoodBin"]
neighborhood_bin_test = pd.DataFrame(index = test_df.index)
neighborhood_bin_test["NeighborhoodBin"] = test_df_munged["NeighborhoodBin"]

#####

numeric_features = train_df_munged.dtypes[train_df_munged.dtypes != "object"].index

# Transform the skewed numeric features by taking log(feature + 1).
# This will make the features more normal.
from scipy.stats import skew

skewed = train_df_munged[numeric_features].apply(lambda x:
skew(x.dropna().astype(float)))
skewed = skewed[skewed > 0.75]
skewed = skewed.index

train_df_munged[skewed] = np.log1p(train_df_munged[skewed])
test_df_munged[skewed] = np.log1p(test_df_munged[skewed])

# Additional processing: scale the data.
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaler.fit(train_df_munged[numeric_features])

scaled = scaler.transform(train_df_munged[numeric_features])
for i, col in enumerate(numeric_features):
    train_df_munged[col] = scaled[:, i]

scaled = scaler.transform(test_df_munged[numeric_features])
for i, col in enumerate(numeric_features):
    test_df_munged[col] = scaled[:, i]

#####

# Convert categorical features using one-hot encoding.
def onehot(onehot_df, df, column_name, fill_na, drop_name):
    onehot_df[column_name] = df[column_name]
    if fill_na is not None:
        onehot_df[column_name].fillna(fill_na, inplace=True)

```

```

dummies = pd.get_dummies(onehot_df[column_name], prefix="_" + column_name)

# Dropping one of the columns actually made the results slightly worse.
# if drop_name is not None:
#     dummies.drop(["_" + column_name + "_" + drop_name], axis=1, inplace=True)

onehot_df = onehot_df.join(dummies)
onehot_df = onehot_df.drop([column_name], axis=1)
return onehot_df

def munge_onehot(df):
    onehot_df = pd.DataFrame(index = df.index)

    onehot_df = onehot(onehot_df, df, "MSSubClass", None, "40")
    onehot_df = onehot(onehot_df, df, "MSZoning", "RL", "RH")
    onehot_df = onehot(onehot_df, df, "LotConfig", None, "FR3")
    onehot_df = onehot(onehot_df, df, "Neighborhood", None, "OldTown")
    onehot_df = onehot(onehot_df, df, "Condition1", None, "RRNe")
    onehot_df = onehot(onehot_df, df, "BldgType", None, "2fmCon")
    onehot_df = onehot(onehot_df, df, "HouseStyle", None, "1.5Unf")
    onehot_df = onehot(onehot_df, df, "RoofStyle", None, "Shed")
    onehot_df = onehot(onehot_df, df, "Exterior1st", "VinylSd", "CBlock")
    onehot_df = onehot(onehot_df, df, "Exterior2nd", "VinylSd", "CBlock")
    onehot_df = onehot(onehot_df, df, "Foundation", None, "Wood")
    onehot_df = onehot(onehot_df, df, "SaleType", "WD", "Oth")
    onehot_df = onehot(onehot_df, df, "SaleCondition", "Normal", "AdjLand")

    # Fill in missing MasVnrType for rows that do have a MasVnrArea.
    temp_df = df[["MasVnrType", "MasVnrArea"]].copy()
    idx = (df["MasVnrArea"] != 0) & ((df["MasVnrType"] == "None") |
(df["MasVnrType"].isnull()))
    temp_df.loc[idx, "MasVnrType"] = "BrkFace"
    onehot_df = onehot(onehot_df, temp_df, "MasVnrType", "None", "BrkCmn")

    # Also add the booleans from calc_df as dummy variables.
    onehot_df = onehot(onehot_df, df, "LotShape", None, "IR3")
    onehot_df = onehot(onehot_df, df, "LandContour", None, "Low")
    onehot_df = onehot(onehot_df, df, "LandSlope", None, "Sev")
    onehot_df = onehot(onehot_df, df, "Electrical", "SBrkr", "FuseP")
    onehot_df = onehot(onehot_df, df, "GarageType", "None", "CarPort")
    onehot_df = onehot(onehot_df, df, "PavedDrive", None, "P")
    onehot_df = onehot(onehot_df, df, "MiscFeature", "None", "Othr")

    # Features we can probably ignore (but want to include anyway to see
    # if they make any positive difference).
    # Definitely ignoring Utilities: all records are "AllPub", except for
    # one "NoSeWa" in the train set and 2 NA in the test set.
    onehot_df = onehot(onehot_df, df, "Street", None, "Grvl")
    onehot_df = onehot(onehot_df, df, "Alley", "None", "Grvl")
    onehot_df = onehot(onehot_df, df, "Condition2", None, "PosA")
    onehot_df = onehot(onehot_df, df, "RoofMatl", None, "WdShake")
    onehot_df = onehot(onehot_df, df, "Heating", None, "Wall")

    # I have these as numerical variables too.
    onehot_df = onehot(onehot_df, df, "ExterQual", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "ExterCond", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "BsmtQual", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "BsmtCond", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "HeatingQC", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "KitchenQual", "TA", "Ex")
    onehot_df = onehot(onehot_df, df, "FireplaceQu", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "GarageQual", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "GarageCond", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "PoolQC", "None", "Ex")
    onehot_df = onehot(onehot_df, df, "BsmtExposure", "None", "Gd")
    onehot_df = onehot(onehot_df, df, "BsmtFinType1", "None", "GLQ")

```

```

onehot_df = onehot(onehot_df, df, "BsmtFinType2", "None", "GLQ")
onehot_df = onehot(onehot_df, df, "Functional", "Typ", "Typ")
onehot_df = onehot(onehot_df, df, "GarageFinish", "None", "Fin")
onehot_df = onehot(onehot_df, df, "Fence", "None", "MnPrv")
onehot_df = onehot(onehot_df, df, "MoSold", None, None)

# Divide up the years between 1871 and 2010 in slices of 20 years.
year_map = pd.concat(pd.Series("YearBin" + str(i+1),
index=range(1871+i*20, 1891+i*20)) for i in range(0, 7))

yearbin_df = pd.DataFrame(index = df.index)
yearbin_df["GarageYrBltBin"] = df.GarageYrBlt.map(year_map)
yearbin_df["GarageYrBltBin"].fillna("NoGarage", inplace=True)

yearbin_df["YearBuiltBin"] = df.YearBuilt.map(year_map)
yearbin_df["YearRemodAddBin"] = df.YearRemodAdd.map(year_map)

onehot_df = onehot(onehot_df, yearbin_df, "GarageYrBltBin", None, None)
onehot_df = onehot(onehot_df, yearbin_df, "YearBuiltBin", None, None)
onehot_df = onehot(onehot_df, yearbin_df, "YearRemodAddBin", None, None)

return onehot_df

# Add the one-hot encoded categorical features.
onehot_df = munge_onehot(train_df)
onehot_df = onehot(onehot_df, neighborhood_bin_train, "NeighborhoodBin", None,
None)
train_df_munged = train_df_munged.join(onehot_df)

# These onehot columns are missing in the test data, so drop them from the
# training data or we might overfit on them.
drop_cols = [
    "_Exterior1st_ImStucc", "_Exterior1st_Stone",
    "_Exterior2nd_Other", "_HouseStyle_2.5Fin",

    "_RoofMatl_Membran", "_RoofMatl_Metal", "_RoofMatl_Roll",
    "_Condition2_RRAe", "_Condition2_RRAn", "_Condition2_RRNn",
    "_Heating_Floor", "_Heating_OthW",

    "_Electrical_Mix",
    "_MiscFeature_TenC",
    "_GarageQual_Ex", "_PoolQC_Fa"
]
train_df_munged.drop(drop_cols, axis=1, inplace=True)

onehot_df = munge_onehot(test_df)
onehot_df = onehot(onehot_df, neighborhood_bin_test, "NeighborhoodBin", None, None)
test_df_munged = test_df_munged.join(onehot_df)

# This column is missing in the training data. There is only one example with
# this value in the test set. So just drop it.
test_df_munged.drop(["_MSSubClass_150"], axis=1, inplace=True)

# Drop these columns. They are either not very helpful or they cause overfitting.
drop_cols = [
    "_Condition2_PosN", # only two are not zero
    "_MSZoning_C(all)",
    "_MSSubClass_160",
]
train_df_munged.drop(drop_cols, axis=1, inplace=True)
test_df_munged.drop(drop_cols, axis=1, inplace=True)

#####

# We take the log here because the error metric is between the log of the
# SalePrice and the log of the predicted price. That does mean we need to

```

```

# exp() the prediction to get an actual sale price.
label_df = pd.DataFrame(index = train_df_munged.index, columns=["SalePrice"])
label_df["SalePrice"] = np.log(train_df["SalePrice"])

print("Training set size:", train_df_munged.shape)
print("Test set size:", test_df_munged.shape)

#####

# XGBoost -- I did some "manual" cross-validation here but should really find
# these hyperparameters using CV. ;-)

import xgboost as xgb

regr = xgb.XGBRegressor(
    colsample_bytree=0.2,
    gamma=0.0,
    learning_rate=0.01,
    max_depth=4,
    min_child_weight=1.5,

n_estimators=7200,
    reg_alpha=0.9,
    reg_lambda=0.6,
    subsample=0.2,
    seed=42,
    silent=1)

regr.fit(train_df_munged, label_df)

# Run prediction on training set to get a rough idea of how well it does.
y_pred = regr.predict(train_df_munged)
y_test = label_df
print("XGBoost score on training set: ", rmse(y_test, y_pred))

# Run prediction on the Kaggle test set.
y_pred_xgb = regr.predict(test_df_munged)

#####

from sklearn.linear_model import Lasso

# I found this best alpha through cross-validation.
best_alpha = 0.00099

regr = Lasso(alpha=best_alpha, max_iter=50000)
regr.fit(train_df_munged, label_df)

# Run prediction on training set to get a rough idea of how well it does.
y_pred = regr.predict(train_df_munged)
y_test = label_df
print("Lasso score on training set: ", rmse(y_test, y_pred))

# Run prediction on the Kaggle test set.
y_pred_lasso = regr.predict(test_df_munged)

#####

# Blend the results of the two regressors and save the prediction to a CSV file.

y_pred = (y_pred_xgb + y_pred_lasso) / 2
y_pred = np.exp(y_pred)

pred_df = pd.DataFrame(y_pred, index=test_df["Id"], columns=["SalePrice"])
pred_df.to_csv('output.csv', header=True, index_label='Id')

```

.







.





.

