

pythonを使った総合的なデータ探索

人生で一番難しいのは自分自身を知ることだ by タレス(ギリシャ時代の哲学者、数学者、宇宙学者、最初の哲学者とよばれる)

データを知るのはデータサイエンスで最も難しいとは言わないけど時間がかかる。だからデータをよく見もせずに分析を始めちゃうんだ。

だから僕はまずやる前にデータ分析のやり方を学ぼうとした。この本の「データを調べる」の章にそって、完璧じゃないけどざっくりとしたデータ分析にのめりこんだ。だからこのカーネルをやったからといってすぐに一人前の研究ができるとは思っていないけど、はじめての人たちにとって役に立つと思うから僕がどうやってデータ分析しているかを共有するね。

変なタイトルつけちゃったけど、このカーネルでやることはこんな感じだ。

1. 問題を理解する それぞれの変数を見て、この問題の意味や重要さについて哲学的に分析しよう
2. 単一変数の研究 従属変数である '販売価格'にフォーカスして、ちょっと詳しくなるう。
3. 多変数の研究 従属変数や独立変数がどうやって関係し合うか理解しよう。
4. 基本のクリーニング 前処理するよ。データの欠損や異常値、カテゴリ変数をうまく対処しよう。
5. 仮説のテスト 多変量解析の技術をつかってデータを仮説にぶっこんで確認してみよう

じゃあ、楽しんでいこう！

In [2]:

```
1 # kaggleパーティに君を招待するね
2 # データ分析に必要なものをimportする
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 from scipy.stats import norm # 正規分布モジュール
8 from sklearn.preprocessing import StandardScaler # データを標準化するもの
9 from scipy import stats
10 # 警告を制御する
11 import warnings
12 warnings.filterwarnings('ignore')
13 %matplotlib inline
```

In [3]:

```
1 # 6缶ケースを持ってくる
2 # データをジュースかお酒扱い?
3 df_train = pd.read_csv('./competitions/house-prices-advanced-regression-techniques/train.csv')
```

In [4]:

```
1 # 項目名がぱっと出せるよ
2 df_train.columns
```

Out[4]:

```
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.で？何すればいいの？

データを理解するには、変数を一つ一つ見て、意味や問題との関係性を理解するんだ。時間がかかるかもしれないけど、この作業でデータの感じを掴めるんだ。（おろそかにしちゃだめだよ！）

分析を規律をもたらすためにエクセルのワークシートに下記の項目を用意してほしい。

- 変数名
- 型 数値型かカテゴリー型 数値は身長や体重みたいに数字で表せるもの、カテゴリーデータは血液型とか男性女性とかの分類のこと。
- 領域 データのセグメント（分類分け）。今回は「建物」、「スペース」、「場所」の3つに分類するよ。「建物」は家の物理的特性に関するもの（鉄筋or木造とか築年数）、「スペース」は家の内容に関するもの（ベットルームの数とか）、場所は家の立地に関するもの
- 予想: 販売価格への影響度合い。「High」、「Middle」、「Low」の3つに分ける。
- 結論: 最終的にどの変数が重要か。（家の売値に重要な要素はどれ？）あとでぱっと見わかるように。expectationと同じ分け方。
- コメント: コメント欄

型とセグメントは見ればわかるけど、大事なのはexpectationでこれが分かれば家の売値の予測がパッとできちゃうんだ。（知らない人からすれば第六感のようにね）それを埋めるために一つ一つ注意深く読んで下記のような質問を自分自身に聞くんだ。

- 自分だったら家を買うときにその変数について考えるか？（たとえば家を建てる時と夢見たときに家の材質が石工単板タイプかなんて考えるかい？（僕はどうでもいいともうけど）
- もし考えたとしたら、それはどれだけ大事でしょうか？()
- 他の変数で説明された情報はないか？（平屋ならスロープがあるかなんて考えなくてもいいでしょ？）

この地道な調査をみっちりやってシートを埋めたら、「予測」が"High"の部分に絞って、販売価格と散布図を作ってみましょう。これで予想が正しかったかを「結論」に埋めていこう！

このプロセスをへて、この問題で下記の変数が重要と結論づけたよ。

- 全体的な素材と仕上りの質
- 建てられた年
- 地下部屋の敷地面積
- 地上のリビングの敷地面積

結論として建物の特性として上の2変数、家の内容として下の2変数が大事だと分かったよ。これはちょっと予想外でした。不動産業界の金言「大事なのは1に立地、2に立地、3に立地」とは真逆だからね。このお手軽な調査はカテゴリーデータに対してちょっと厳しいかもね。たとえば僕はご近所は大事だと思ったんだけど分析したら大事じゃないって分かったよ。これはカテゴリーデータの可視化に適した箱ひげ図じゃなくて散布図を使ったからだね。データの可視化の仕方しだいでは結果に影響するんだよ。

でもこのレッスンで大事なことは僕たちはデータやその関連度について少し立ち止まって考えたってこと、だからゴールにたどり着いんだ。おしゃべりは終わりにして、早速やってみよう！

2.最初の最初？ 販売価格について分析してみよう

販売価格は分析の目的だ。パーティに出かけるときのようだ。パーティへいく理由は大抵"女性"だ。("女性"の部分は自分の好みに置き換えてくれ)

女性のたとえをつかって、ちょっとした小話を作ってみよう。題して「どうやって販売価格に出会うか」。

kaggleパーティが始まって、ダンスパートナーを探している。ダンスフロアを探し回っていて、バーの近くにダンスシューズを履いた少女に出会った。我々はモデル予測やデータ分析コンペをたくさんやってきたいわばプロだ。少女に話しかけるのなんてわけない。なん

「やあ、ぼくkaggler。あなたは？「販売価格」さん？いい名前だね。もっと情報をくれないかな？ぼくは誰と誰が仲良くなるかを確率的に計算するモデルを作りたいんだ。そしてそれを応用したいんだ。

In [5]:

```
1 # describe()関数で「販売価格」さんの情報が見れるよ！
2 #descriptive statistics summary
3 # その項目の値についての概要情報を見れる
4 # count :データの個数
5 # mean: 平均値
6 # std: 標準偏差
7 # min: 最小値
8 # 25% 50% 75%: 下から1/4,1/2,3/4のところにある数値
9 # max: 最大値
10 # Name: 項目名 dtype:型
11 df_train['SalePrice'].describe()
```

Out[5]:

```
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
```

In [6]:

```
1 df_train.head() # 先頭から5件のデータを表示する。データをざっと見ることができる基本コマンド
```

Out[6]:

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

5 rows × 81 columns

In []:

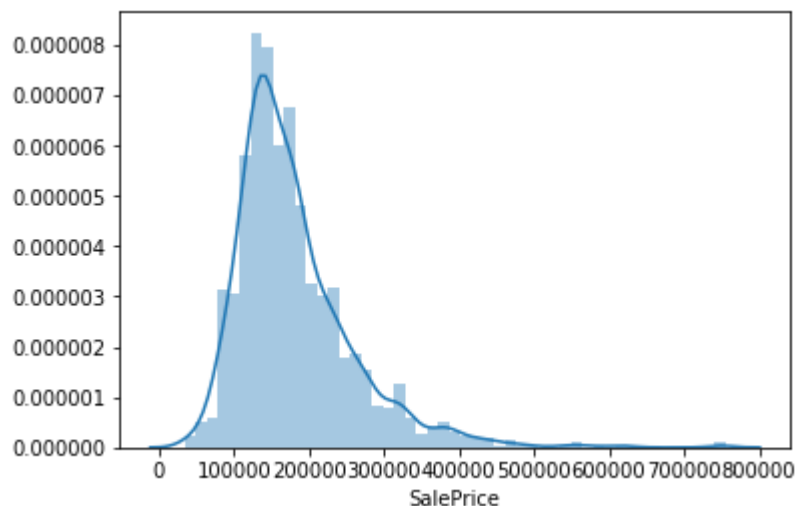
```
1 なるほど、最安値はゼロより大きいね。すばらしい、あなたは私のモデルをぶっ壊すかもしれない性質の一つを持って
2  なにか僕に見せられる写真はない？わからないけど、ビーチにいたときとか事務の写真とかさー。
```

In [29]:

```
1 # ヒストグラム 下記に引数のデフォルト値や解説あり
2 # https://seaborn.pydata.org/generated/seaborn.distplot.html
3 # ヒストグラムhist=Trueと確率密度関数kde=Trueとラグプロット(どこに値があるか軸上にぶろっとする)=False)
4 sns.distplot(df_train['SalePrice'])
```

Out[29]:

<matplotlib.axes._subplots.AxesSubplot at 0x1a1e98f198>



ああ、君がちょっと席をはずしているときにいろいろのぞいちゃったんだ、すごいね！君は

- 正規分布からはずれているね
- 正の歪度だから左寄りだね
- 山のてっぺんでピークがわかるね

これは面白い。「販売価格」さん、もっといろいろ調べさせてよ！

In [8]:

```

1 #skewness and kurtosis
2 # 歪度 skewness 正規分布からどれだけ歪んでいるか 左寄りのとき正の値をとり、右寄りのとき負の値をtoru
3 # 尖度 kurtosis とがっているか 尖るほど値が大きくなる
4 # SalesPriceは左寄りで尖っていることがわかる
5 print("Skewness: %f" % df_train['SalePrice'].skew())
6 print("Kurtosis: %f" % df_train['SalePrice'].kurt())

```

Skewness: 1.882876

Kurtosis: 6.536282

すごいや！もし僕の計算が正しいければ成功率は97%だよ。きみとはまた会いたいよ。これぼくの電話番号だから覚えていてよ。もし金曜暇ならかけてきてね！じゃあまたね！※なんて自分勝手なkaggle野郎だ。

販売価格さんの仲間たちと興味

地形を選ぶのは戦術的知恵だ。まもなく「販売価格」さんが君の眼前に現れるだろう。心してかかれよ。追っ手はいない。こわがらなくても個別の変数をガンガン調査するだけさ！我々の研究を最大限活用するために、私たちは共通の友人のプロフィールに注意深く目を通し、共通の興味に焦点をあてます。

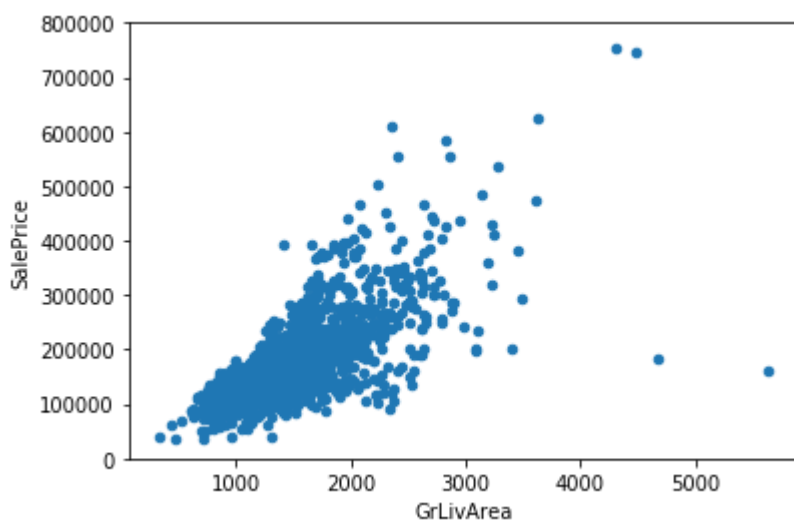
数値データ間の関係

In [6]:

```

1 # grlivareaとsalepriceの関係を散布図にする
2 var = 'GrLivArea'
3 # 項目を横並びにならべる。axis=0だと縦に並べる
4 data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
5 # yの範囲を0,800000にして散布図を表示
6 data.plot.scatter(x=var, y='SalePrice', ylim=(0, 800000));

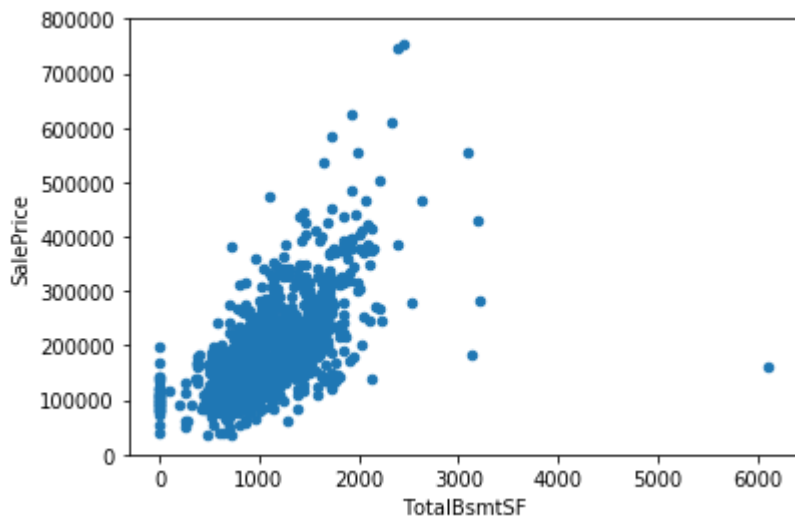
```



うーん、「販売価格」と「地上リビングの敷地面積」さんは線形関係をもつ古くからの友人のようだ。「合計地下の敷地面積」君はどうだろう？

In [10]:

```
1 # totalbsmtsf/salepriceを散布図で表す
2 var = 'TotalBsmtSF'
3 data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
4 data.plot.scatter(x=var, y='SalePrice', ylim=(0, 800000));
```



「TotalBsmtSF」君も「販売価格」さんとお友達みたい。しかしこれはもっと感情的な関係ね。指数関数的な感じね。さらにときどき全然相関がなくなるみたいね。

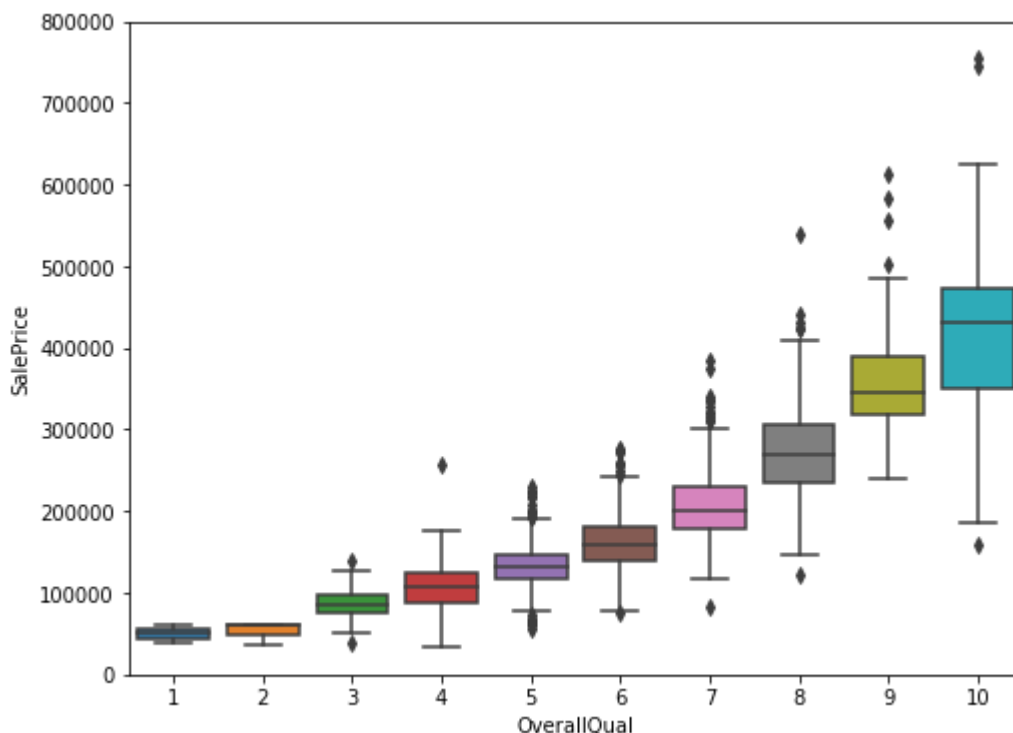
カテゴリー的な特徴の関係

In [11]:

```

1 #ヒゲ箱プロット overallqualと販売価格
2 # 箱ひげ図の見方
3 # それぞれのカテゴリデータをまとめて
4 # 上ヒゲが最大値、下ヒゲが最小値、箱の上が第三四分位数、真ん中が第二四分位数
5 # 下が第一四分位数
6 # 黒豆さんは何?(@@@調査中です@@@)
7 var = 'OverallQual'
8 data= pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
9 f, ax = plt.subplots(figsize=(8, 6)) # 画面サイズを8:6 横:縦で表示
10 fig = sns.boxplot(x=var, y='SalePrice', data=data)
11 fig.axis(ymin=0, ymax=800000);
12
13 # ※f, ax = の部分は figure axisオブジェクトをアンパック代入
14 # たとえばfig.savefig('yourfilename.png')で保存できたり
15 # するので慣例的にこう書いておく(axの使い方は調べきれていない)
16 # 箱の上下にあるのは外れ値箱の1.5倍以上の大きさの点は外れ値とされる
17 # whisパラメータで調節でき whis='range'とすると外れ値にしない

```



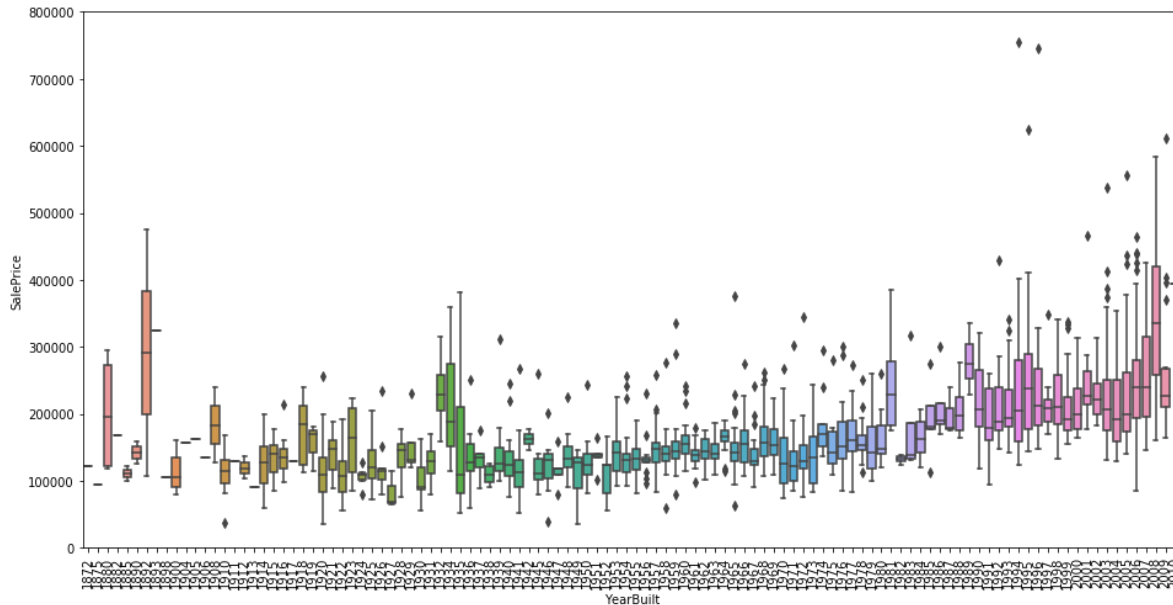
かわいい少女のように「販売価格」さんは「全体的な素材や仕上がりの質」を楽しんでいる。メモ：最初のデータにマクドナルドはどうか？（だめでしょ？）

In [34]:

```

1 var = 'YearBuilt'
2 data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
3 f, ax = plt.subplots(figsize=(16,8)) #描画サイズのデフォルトは8:6
4 fig = sns.boxplot(x=var, y="SalePrice", data=data)
5 fig.axis(ymin=0, ymax=800000);
6 plt.xticks(rotation=90); # x軸に関する指定 ここでは年数が横に並べると重なって読めないなので90度回転させ

```



そこまで強い傾向ではないですが、古い家より新しい家の方が高値で売られていますね。

注意：販売価格が実質価格か分かっていません。実質価格はインフレの影響を取り除きます。もし実質価格になっていない場合はそうすべきで、それによって時系列での価格を正しく扱えます。

まとめ

話がそれましたが、結論づけると、

- 「地上のリビングの敷地面積」と「地下の敷地面積」は販売価格と比例関係にありそうです。正の相関があって、片方の変数の値が増えともう片方も増えるといった具合です。TotalBsmtSFの場合は値が大きくなるとより販売価格が高くなっていくようですね。
- 「全体的な素材と仕上りの質」と「建てられた年」もどちらも販売価格と関係があるようです。「全体的な素材と仕上りの質」の場合は特に強い関係性があり、箱ひげ図が販売価格が全体の質にそって増えることを示しています。

今まで4つ変数しか分析してこなかったけど、他の変数についてもやるべき分析がたくさんあります。肝は正しい変数を選び出すこと（特徴選択）で、難しい関係性を定義すること（特徴量エンジニアリング）じゃないよ。

つまりたくさんの特徴から必要な特徴を選り分けることだ！

3. 落ち着いてスマートにやろう！

これまでは直感にしたがって重要だと思った変数を分析してきました。客観的に分析しているつもりでも、出発点は主観的だったといわざるをえません。このアプローチはエンジニアとして癪に触るよね。私から言えるのは規律を守って、こうなんじゃないかって主観を働かせないようにすることだ。それには訳があります。構

造工学では主観的にやろうとすると、法則がガラガラと崩れ去り痛い目を見るでしょう。

だから主観的にになりたいという慣性に打ち勝ってぜひ客観的な分析を行うようにしてください。

プラズマスープ

宇宙の始まりの時、そこはプラズマスープだけだった。宇宙を研究し始めたときはその事実は推測にすぎなかった。しかし、今日までの研究で科学は宇宙で何が起きたかを鮮やかに描ききっている。

宇宙を探索するには、プラズマスープを意味する実践的なレシピから始めるのがいいでしょう。

- 相関行列(ヒートマップ) *個々の値のデータ行列を色として表現した可視化グラフ
- 販売価格の相関行列 (ズームしたヒートマップ)
- もっとも相関の高いものを散布図でプロットしろ (ジャガースタイルに移行せよ)

相関行列(ヒートマップ)

相関行列

それぞれの項目ごとの相関係数をもとに色付けした図。

相関係数

相関係数とは2種類のデータの関係を示す指標で、-1から1のあたいで示され、

- -1か1に近づくほど強い相関がある。逆に0はほぼ相関がない。
- 正の値のときは正の相関があるといい、どちらかの値が増加するともう一方の値が増加傾向
- 負の値のときは負の相関があるといい、どちらかの値が増加するともう一方の値が減少する傾向

相関係数 r は S_{xy} 共分散、 S_x x項目の標準偏差、 S_y y項目の標準偏差として

$r = \frac{S_{xy}}{S_x S_y}$ となる。

次に標準偏差は割愛して共分散について説明する。

共分散

2組の対応するデータ間の関係を表す数値である。データのペア $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ があつたとして、

共分散 S_{xy} は

$S_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \overline{x})(y_i - \overline{y})$ で求められる。

共分散は `np.cov()` で求められるが、中身の計算を理解するために一つ一つ計算してみる。

In [21]:

```
1 import math
2 # 試しにSalePriceとGrLivAreaの共分散と相関係数を求めている
3 col_SP = df_train['SalePrice']
4 col_GLA = df_train['GrLivArea']
5
6 mean_SP = np.mean(col_SP)
7 mean_GLA = np.mean(col_GLA)
8
9 std_SP = math.sqrt(sum([(n - mean_SP)**2 for n in col_SP])/len(col_SP))
10 # 計算した値とnp.stdが正しいか確認
11 print(std_SP, np.std(col_SP))
12 std_GLA = np.std(col_GLA)
13
14 cov = 0
15 for x,y in zip(col_SP, col_GLA):
16     cov += (x-mean_SP)*(y-mean_GLA)
17
18 cov/= len(col_SP)
19
20 corr = cov / (std_SP*std_GLA)
21
22 print("共分散は{}, 相関係数は{}".format(cov,corr))
```

79415.29188606751 79415.29188606751

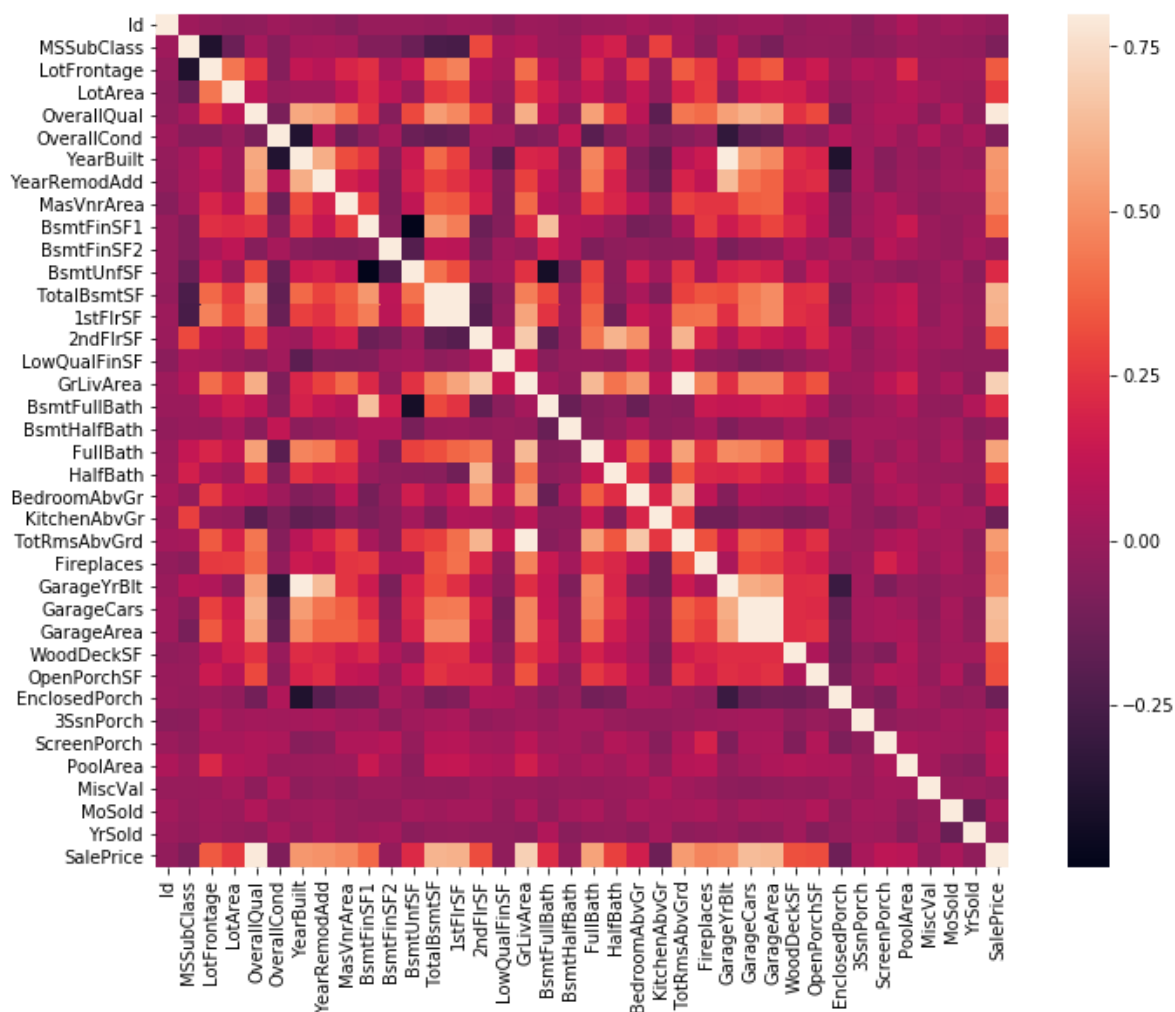
共分散は29561605.190672714、相関係数は0.7086244776126512

In [14]:

```

1 # 相関行列
2 # 一般的に0.7以上でかなり強い正の相関があるとされる 白っぽいオレンジあたりから?
3 corrmat = df_train.corr()
4 f, ax = plt.subplots(figsize=(12, 9))
5 sns.heatmap(corrmat, vmax=.8, square=True);
6 # vmax=0.8にすると0.8より大きい値の色が固定になる。これでかなり強い相関はひと目。
7 # squareで各ブロックを正方形にする

```



私見だが、ヒートマップは全体やその関係をパッと俯瞰するのに最適なやり方だ。

一目で目を引くものが2つあります。'TotalBsmtSF'と'1stFlrSF'、'GarageX'系の変数だ。これらは重要な相関があるのがわかります。多重共線性（マルチコ）だとわかります。これは同じ情報だとわかるので、それを削ったりするのに使える便利なものです。

<http://xica.net/vno4ul5p/> (<http://xica.net/vno4ul5p/>) マルチコについては上記参照。同じような項目を入れると起きる現象で基本的に項目を絞ることで対策できる。

他に販売価格の相関関係を見てみましょう。今までよく見てきた「GrLivArea」「TotalBsmtSF」「OverallQual」が高い相関があることが見てとれます。ほかにも取り上げるべき変数がいくつかあるので、それを次にやってみましょう。

In [15]:

```
1 col_check = corrmat.nlargest(10, 'SalePrice')['SalePrice'].index
```

In [16]:

```
1 # NumPyのcorrcoefは行と行の相関係数を計算するため転置する必要あり
2 np.corrcoef(df_train[col_check].values.T)
```

Out[16]:

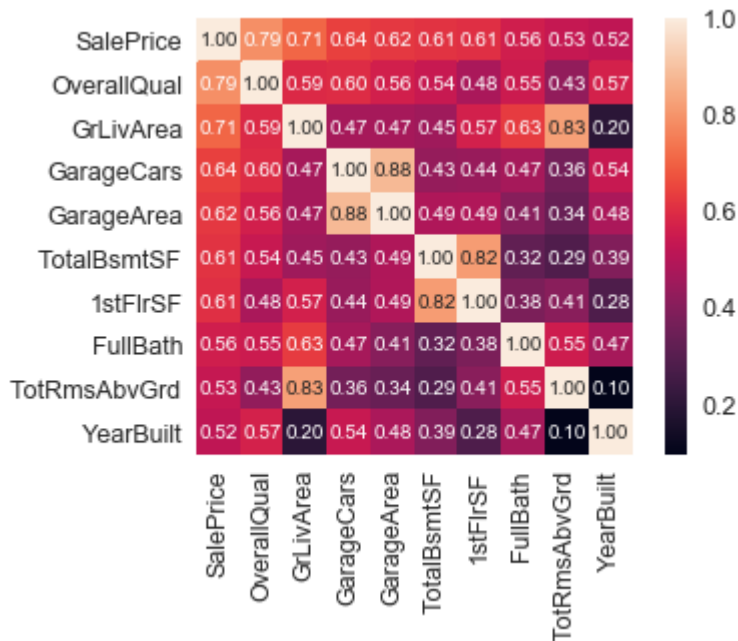
```
array([[1.          , 0.7909816 , 0.70862448, 0.6404092 , 0.62343144,
        0.61358055, 0.60585218, 0.56066376, 0.53372316, 0.52289733],
       [0.7909816 , 1.          , 0.59300743, 0.60067072, 0.56202176,
        0.5378085 , 0.47622383, 0.55059971, 0.42745234, 0.57232277],
       [0.70862448, 0.59300743, 1.          , 0.46724742, 0.46899748,
        0.4548682 , 0.56602397, 0.63001165, 0.82548937, 0.19900971],
       [0.6404092 , 0.60067072, 0.46724742, 1.          , 0.88247541,
        0.43458483, 0.43931681, 0.46967204, 0.36228857, 0.53785009],
       [0.62343144, 0.56202176, 0.46899748, 0.88247541, 1.          ,
        0.48666546, 0.48978165, 0.40565621, 0.33782212, 0.47895382],
       [0.61358055, 0.5378085 , 0.4548682 , 0.43458483, 0.48666546,
        1.          , 0.81952998, 0.32372241, 0.28557256, 0.391452  ],
       [0.60585218, 0.47622383, 0.56602397, 0.43931681, 0.48978165,
        0.81952998, 1.          , 0.38063749, 0.40951598, 0.28198586],
       [0.56066376, 0.55059971, 0.63001165, 0.46967204, 0.40565621,
        0.32372241, 0.38063749, 1.          , 0.55478425, 0.46827079],
       [0.53372316, 0.42745234, 0.82548937, 0.36228857, 0.33782212,
        0.28557256, 0.40951598, 0.55478425, 1.          , 0.09558913],
       [0.52289733, 0.57232277, 0.19900971, 0.53785009, 0.47895382,
        0.391452  , 0.28198586, 0.46827079, 0.09558913, 1.          ]])
```

In [18]:

```

1 # 販売価格との相関行列(ヒートマップ)
2 k=10 # 変数の数
3 # 関数の意味を記述する nlargest 大きい順にいくつかとる <-> nsmallest
4 cols = corrmatrix.nlargest(k, 'SalePrice')['SalePrice'].index
5 # 販売価格との相関係数が大きい順に10に並べて、そのインデックスをcolsに代入
6 cm = np.corrcoef(df_train[cols].values.T) # 販売価格との相関係数が大きい順10項目の値で再度相関係数
7 # なぜ転置をとるのか http://lofas.hatenablog.com/entry/2015/02/09/150552
8 sns.set(font_scale=1.25) # 項目の名前のフォントを1.25倍する
9 hm = sns.heatmap(cm, cbar=True, annot=True, square=True, fmt='.2f', annot_kws={'size': 10}, yti

```



見てみると、関係の強い変数が現れる。それぞれについて私が思ったことは、

- 「全体的な素材と仕上がりの質」「家の中の面積」「総敷地面積」は販売価格と強い関係性がある チェックや！
- 「ガレージに駐車できる数」と「ガレージの面積」は「販売価格」と強い関係がある。しかし、ガレージが大きいからたくさん駐車できるだけで、この二つは双子の関係なので関係性がつよい「ガレージに駐車できる数」のみ残す。
- 「総敷地面積」と「1階の面積」も同じようなものなので、最初の直感で選んだ「総敷地面積」を残す。
- 「地上の浴槽・シャワー・洗面台・便器の4点が備え付けられた部屋」？本当？
- 「バスルームを除いた部屋の数」と「家の中の面積」も兄弟みたいなもの。チェルノブイリのデータセットなの？
- 「建てられた年」か。それも相関があるね、もしかして時系列分析しなきゃいけないかもって思ってる。それはあなたへの宿題ね！

散布図に進もう！

販売価格と相関が強いものを散布図にしよう（ジャガースタイルで動こう）

見るべきものは準備できた。最初に散布図を見た時は感動したよ。たくさんの情報がぎっしり詰まっていますばらしかった。ありがとうseaborn！君のお陰でジャガーのように動けるよ。

In [24]:

```

1 # 散布図
2 sns.set()
3 cols = ['SalePrice', 'OverallQual', 'GrLivArea', 'GarageCars', 'TotalBsmtSF', 'FullBath', 'YearBuilt']
4 sns.pairplot(df_train[cols], size=2.5)
5 plt.show();

```



だいたいグラフを見る前から分かっているんだけど、プロットを見るとより妥当な考えが分かったりするんだ。

「敷地の総合面積」と「家の中の面積」が気になったんだけど、比例関係があって上に境界があるみたいだ。地下の面積は地上とだいたい一緒だけど、地上の面積より地下の面積が大きいことはないよ。（バンカーを買う時以外は）「建てられた年」と「販売価格」は恥ずかしがり屋の指数関数のようで上限が徐々に上がっていくね。

ローシャットテストは済んだね。ここからは欠損値を見ていこう！

4. 欠損値

欠損値について考えるときに重要な質問は

- 欠損値はどれくらいあるか？
- 欠損値はランダムなのかパターンがあるのか？

この質問の答えは実用的な理由で重要です。なぜなら欠損値はサンプルサイズが減ることを意味するからです。これはデータ分析にすむのを阻みます。さらに、欠損値が偏っていないか、不都合な真実を隠していないか保証する必要があるんです。

In [30]:

```
1 # Nanデータの数量を集めて項目ごとの欠損値の数をまとめて降順でtotalに格納
2 total = df_train.isnull().sum().sort_values(ascending=False)
3
4 # 欠損率を降順でpercentに格納
5 percent = (df_train.isnull().sum()/df_train.isnull().count()).sort_values(ascending=False)
6
7 # 欠損数と欠損率をmissing_dataとしてそれを上位20表示する
8 missing_data = pd.concat([total, percent], axis=1, keys=['Total', 'Percent'])
9 missing_data.head(20)
```

Out[30]:

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

欠損データをどう扱うかを理解するために分析してみましょう。

データが15%以上欠損しているとき、類似の変数を消してなかったことにしましょう。このケースでは欠損データに何かを埋めるためになにかするといったことはない。これによっていくつかデータを消せる。これはデータを失ったということ？違うと思います。これらの変数はあまり重要ではないってこと。家を買う時にあまり考えないからデータが集まらなかったと考えてもいいかもね。さらに注意深く見ていくといくつかの変数は外れ値だと思われるので喜んではずしましょう！

のこりのケースではガレージに関する変数は同じくらい欠損データがありますね。だけど5%の欠損率でガタガタ言うよりはガレージの駐車できる車の数だけみれば事足りますね。同じことが地下に関する変数にも当てはまります。

石造りのレンガの面積やタイプに関しては重要じゃないと考える。さらにこれらはすでに考慮した「建てられた年」と「総合的な質」と強い関係がある。故にこの変数を削っても情報は失われません。

まとめると、欠損データを扱うために、「電気システム」を除いたすべての欠損データを削除します。電気システムは外れ値もるとも観測自体をやめよう。(??)

In [20]:

```
1 df_train.Electrical.value_counts()
```

Out[20]:

```
SBrkr    1334
FuseA      94
FuseF     27
FuseP      3
Mix         1
Name: Electrical, dtype: int64
```

In [31]:

```
1 df_train = df_train.drop((missing_data[missing_data['Total']> 1]).index, 1)
2 df_train = df_train.drop(df_train.loc[df_train['Electrical'].isnull()].index)
3 df_train.isnull().sum().max() # ちゃんと欠損値がなくなったか確認
```

Out[31]:

0

うそつき！（外れ値？）

外れ値も注意すべきものだ。なぜか？外れ値はモデルにかなりの影響を与えうるし、特定の振る舞いについての知見を与えてくれる重要な情報源にもなりうるからね。

外れ値は複雑な科目で注意深く見る必要がある。ここで「販売価格」の標準偏差と散布図を通してサッと見ていきましょう。

単一変数の分析

大事なのは外れ値を探すために閾値を設定することです。それをするのに、データを標準化しましょう。ここでいうデータの標準化とは、データを平均0、標準偏差1とすることです。

データの標準化とは

データセットを平均0、標準偏差1にすること

難しいことはなく 各データから平均を引いて、標準偏差で割る だけ。

In [34]:

```

1  # データの標準化
2
3  # 関数の中身を記述する
4  saleprice_scaled = StandardScaler().fit_transform(df_train['SalePrice'][:,np.newaxis]);
5  # np.newaxisはNoneの参照なので[:,None]でも同じだが軸を追加する意味になり、コードの分かりやすさの面から
6  low_range = saleprice_scaled[saleprice_scaled[:,0].argsort()][:10]
7  high_range = saleprice_scaled[saleprice_scaled[:,0].argsort()][-10:]
8  print('outer range (low) of the distribution:')
9  print(low_range)
10 print('\nouter range (high) of the distribution:')
11 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]]

```

「販売価格」に

- 低い範囲の値は似通っていて、平均からそんなに離れていない - 高い範囲の値はずっと離れていて、実際 7 この値は範囲外です。

さて、これらの値を外れ値として考慮しないようにするが、注意しておくべきだ。??

二変量分析

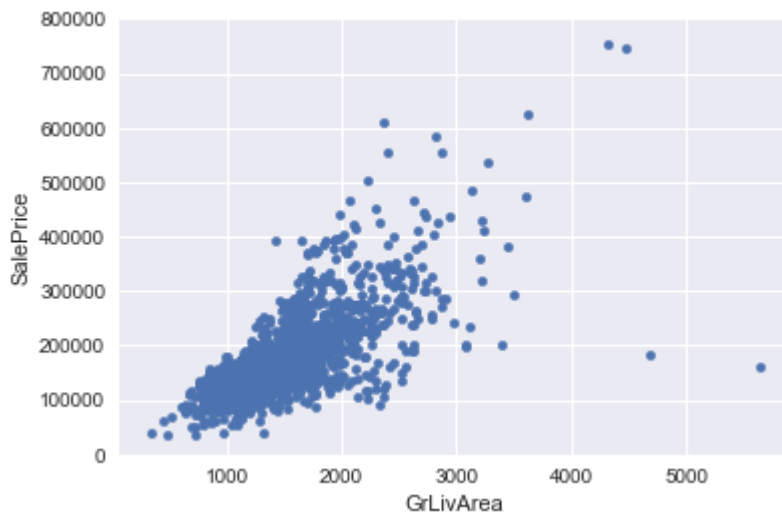
下の散布図はもういやってほど見てきましたね。しかし、新しい観点から見てみると、いつも発見がありますよね。見方を変えるのは80IQほどの価値があるよ。 (by Alan Kay)

In [35]:

```

1 # 二変量分析 販売価格と地上の部屋の面積
2 var = 'GrLivArea'
3 data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
4 data.plot.scatter(x=var, y='SalePrice', ylim=(0,800000));

```



これから分かることは

- 「地上の部屋の面積」が大きくなると次第におかしくなり、群れに従わなくなる。（広がっていく）なぜそうなるのかじっくり考える。もしかして農地に関してそれが安さの原因かも。??わからないが、私はこれらの2つの点が典型的事例を代表するものではないと自信がある。ゆえにそれらを異常値としています。
- 7つの観測点 2つの これらは2つの特例だけど、トレンドにはしたがっているのでキープします。??

In [36]:

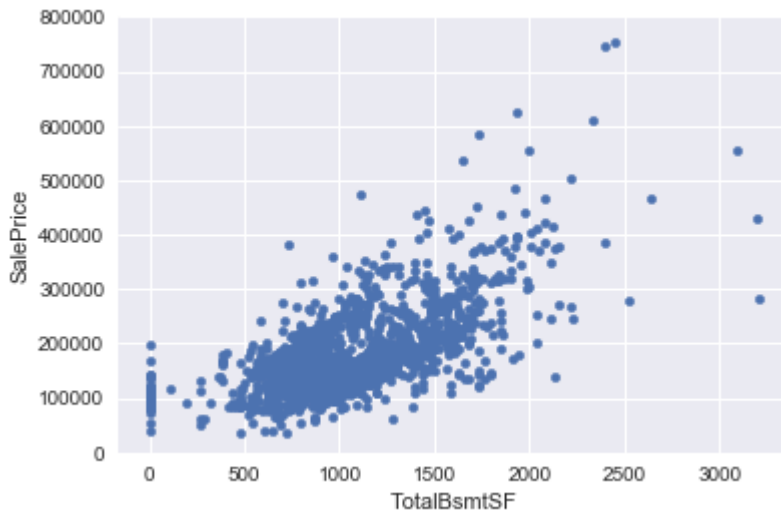
```

1 # 点の削除
2
3 df_train.sort_values(by = 'GrLivArea', ascending=False)[:2]
4 df_train = df_train.drop(df_train[df_train['Id'] == 1299].index)
5 df_train = df_train.drop(df_train[df_train['Id'] == 524].index)

```

In [37]:

```
1 # 2変量分析「販売価格」と「地上の住居面積」
2 var = 'TotalBsmstSF'
3 data = pd.concat([df_train['SalePrice'], df_train[var]], axis=1)
4 data.plot.scatter(x=var, y='SalePrice', ylim=(0, 800000));
```



いくつか観測値を排除したいと思うかもしれませんが、その価値はないと思います。 あっても問題ないので放置します。

5.核心に迫るよ

??は誰だ?って話のほとんど

「販売価格」は誰だ?

質問への答えは多変量解析のための統計的根拠の基礎となる仮定をテストすることで分かるはずだ。 これまでデータを綺麗にして「販売価格」さんについてたくさん発見があった。

ここからさらに深くもぐって、「販売価格」さんに多変量のテクニックを駆使してどう統計的仮説を当てはめるかだ。

- 正常か ここでいう正常は正規分布にそっているかだ。これはいくつかの統計テストが正規分布に基づいているからです。ここでは単変数の正常さをチェックしてきた（これは限られたアプローチです）ただ単変数の正常性は多変数の正常性を保証するとは限らないことを覚えておいてください。けど助けにはなるから。 また200データ以上のサンプルがあるときは正常性はそこまで問題にならない、ただ正常性に関してちゃんとやっておけば多くの問題をさけられそれが分析を行う主要な理由になる。??
- 均質性 均質性は予測変数の範囲に依存変数があるという仮説をいいたい。

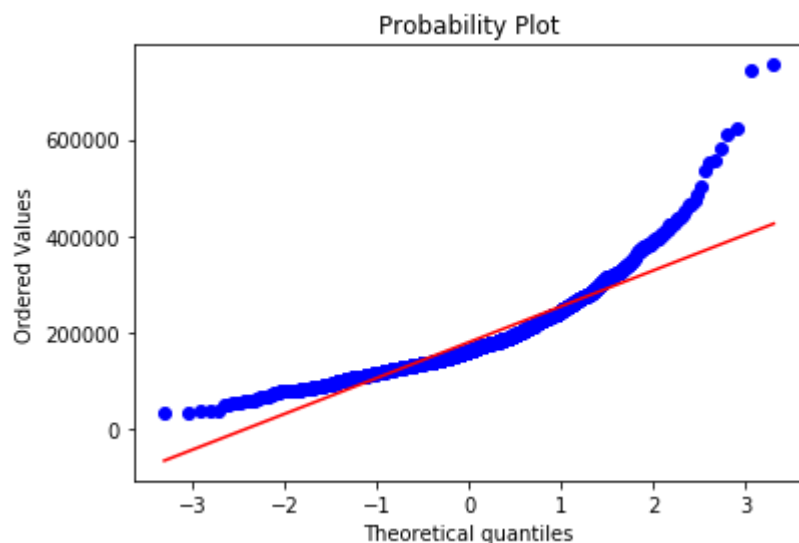
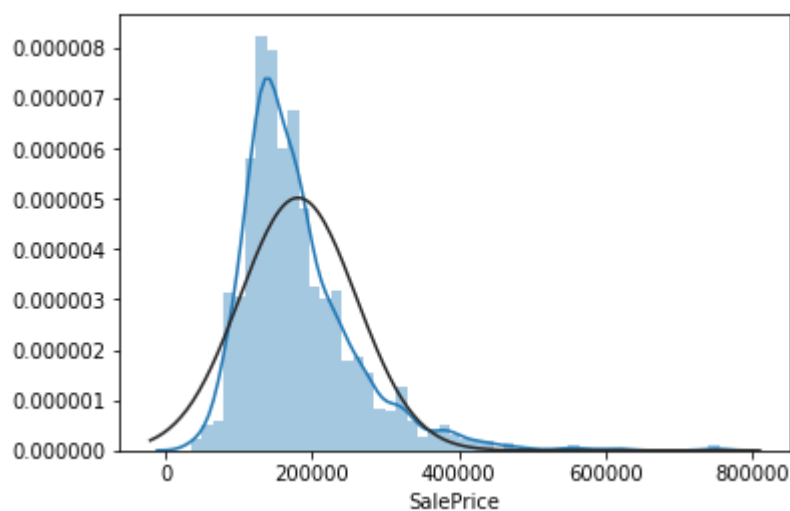
何言っているか分からない?じゃあ行動してみよう!

正常性の調査

- ヒストグラム 歪度と尖度
- 正規確率プロット 正規分布に従っているかを確認できる

In [8]:

```
1 # ヒストグラムと 正規確率プロット
2 sns.distplot(df_train['SalePrice'], fit=norm);
3 fig = plt.figure()
4 res = stats.probplot(df_train['SalePrice'], plot=plt)
```



見てみると「販売価格」は正常でないね。ピークがあって、正の歪度があって、斜めの対角線に従っていないね。

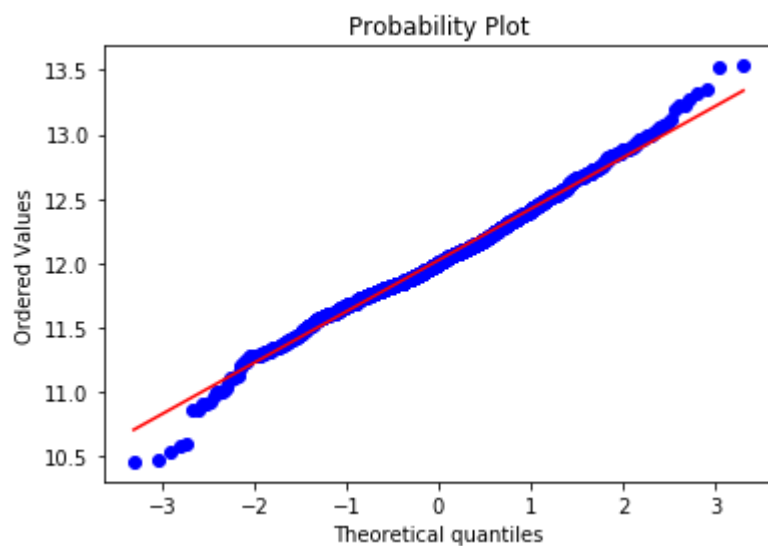
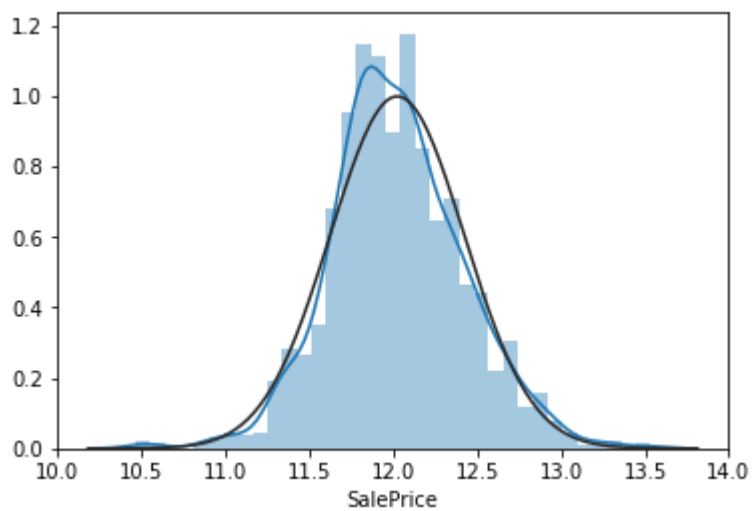
だけどまだ負けたわけじゃない。すごいテクニックがあって正の歪度がある場合は対数をとればうまくいくんだ。すけーだろ！

In [9]:

```
1 # 対数を取る
2 df_train['SalePrice'] = np.log(df_train['SalePrice'])
```

In [10]:

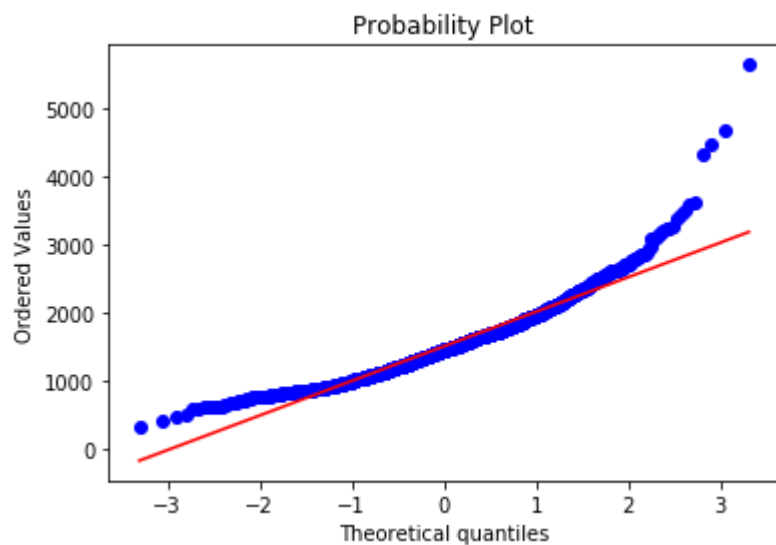
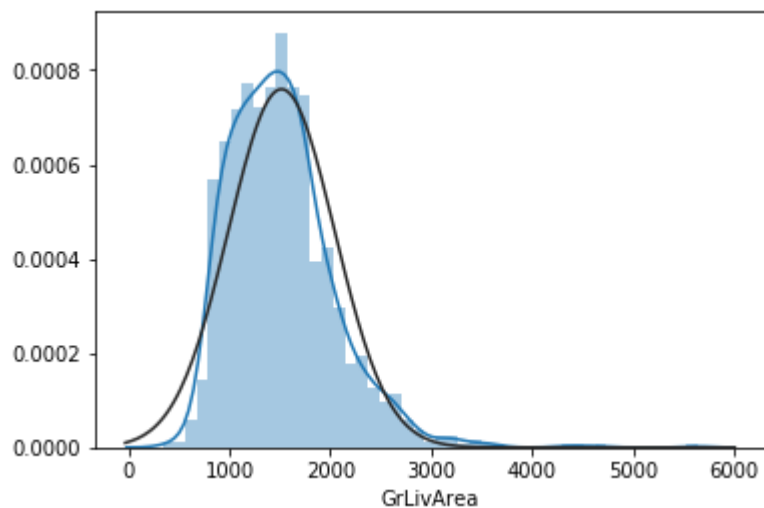
```
1 # どうなったか確認
2 sns.distplot(df_train['SalePrice'], fit=norm);
3 fig = plt.figure()
4 res = stats.probplot(df_train['SalePrice'], plot = plt)
5
```



やったぜ！「地上のリビングの敷地面積」も同じようにやってみよう！

In [11]:

```
1 # まずそのまま見てみる
2 sns.distplot(df_train['GrLivArea'], fit=norm);
3 fig = plt.figure()
4 res = stats.probplot(df_train['GrLivArea'], plot=plt)
```



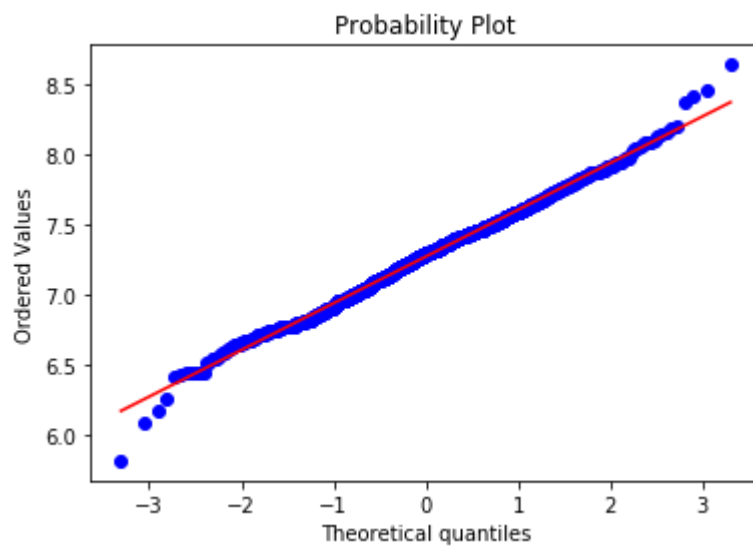
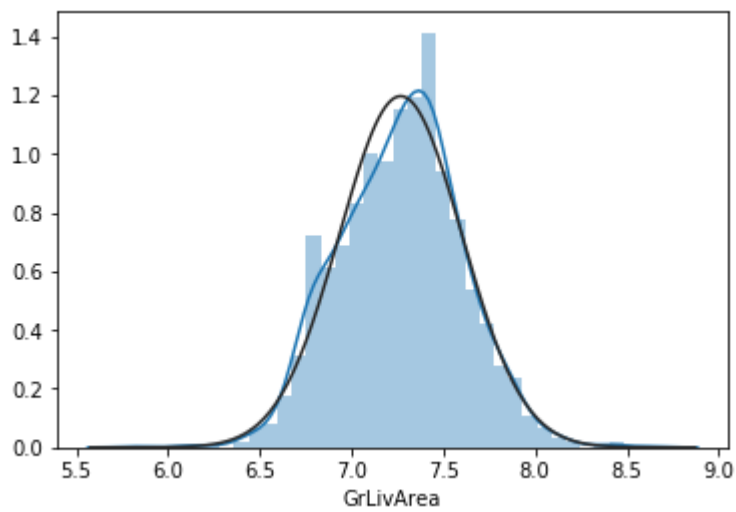
ちょっとゆがんでるな！魔法でちょちょい

In [12]:

```
1 # 対数を取る
2 df_train['GrLivArea'] = np.log(df_train['GrLivArea'])
```


In [14]:

```
1 # 効果を確認
2 sns.distplot(df_train['GrLivArea'], fit=norm)
3 fig = plt.figure()
4 res = stats.probplot(df_train['GrLivArea'], plot=plt)
```



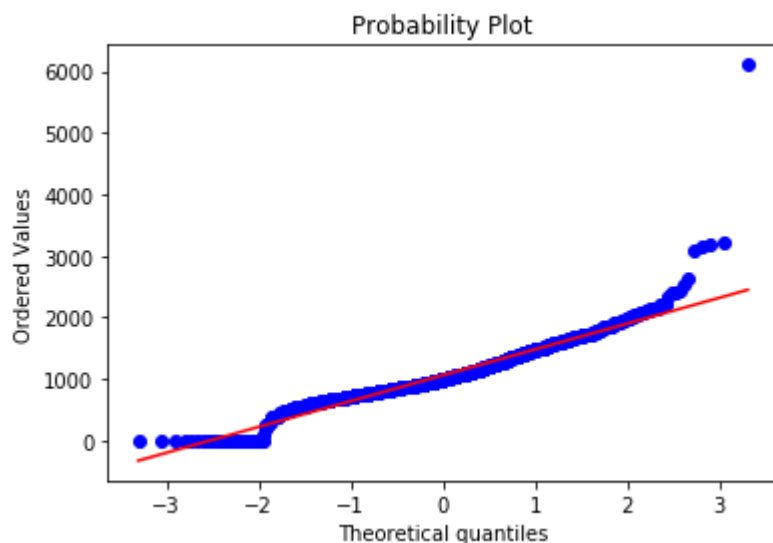
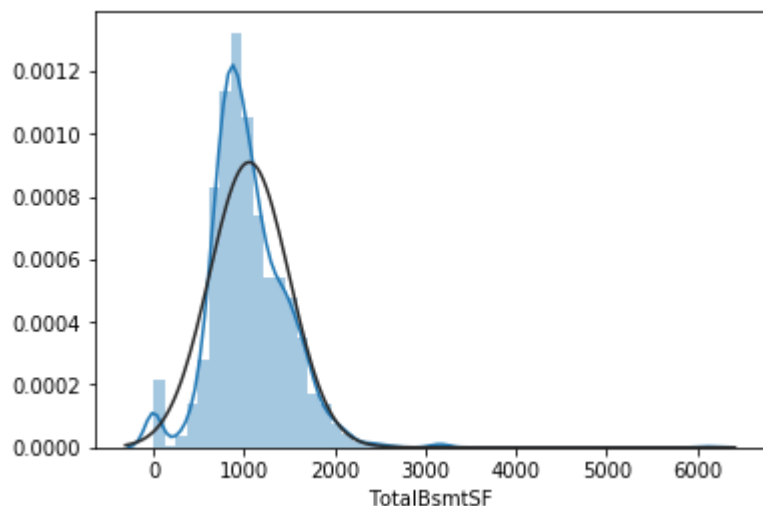
次！

In [15]:

```

1 sns.distplot(df_train['TotalBsmtSF'], fit=norm);
2 fig = plt.figure()
3 res = stats.probplot(df_train['TotalBsmtSF'], plot=plt)

```



なかなか手ごわいやつがきたね！どうすりゃいいかね？

- 歪みが少ない
- 地下がない家もあるから0付近にかなり数がある
- 0は対数をとれない（マイナス無限大になっちゃう）から大問題！

対数をとるために0かそうでないかを区別して、0でない値にだけ対数をとると0であるという情報を損なわない。

私はこのアプローチが正しいか分からない。けどうまく動いてる。これをハイリスクエンジニアリングとよんでいます。

In [16]:

```

1 # 値が0より大きいものを1、それ以外を0とした新しいカテゴリ変数を作る
2 df_train['HasBsmt'] = pd.Series(len(df_train['TotalBsmtSF']), index=df_train.index)
3 df_train['HasBsmt'] = 0
4 df_train.loc[df_train["TotalBsmtSF"]>0, 'HasBsmt'] = 1

```

In [18]:

```

1 # ここで対数を取る
2 df_train.loc[df_train['HasBsmt']==1, 'TotalBsmtSF'] = np.log(df_train['TotalBsmtSF'])

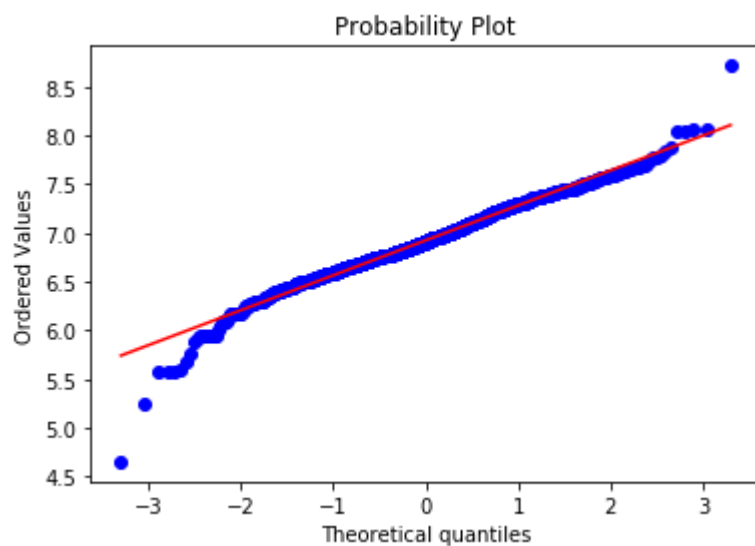
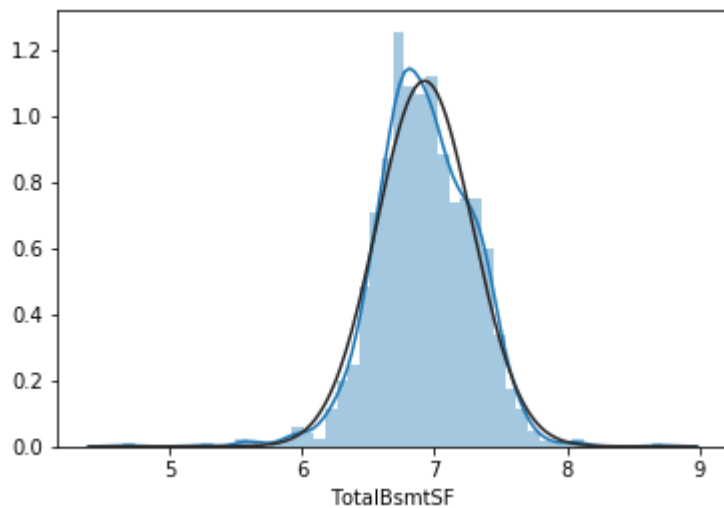
```

In [20]:

```

1 # 確認してみる
2 sns.distplot(df_train[df_train['TotalBsmtSF']>0]['TotalBsmtSF'], fit=norm);
3 fig = plt.figure()
4 res = stats.probplot(df_train[df_train['TotalBsmtSF']>0]['TotalBsmtSF'], plot=plt)

```



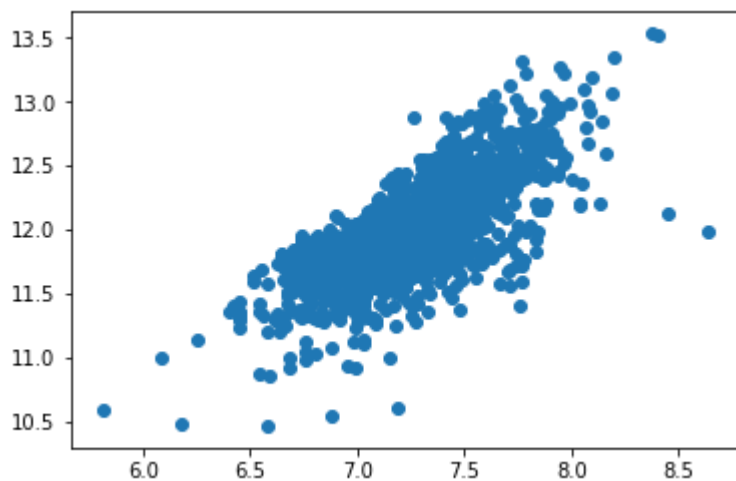
一発で均質性が正しいか調査

2つのメトリック変数の均質性を確認するアプローチはグラフにすること。

例によって「販売価格」と「地上のリビングの敷地面積」を見てみよう

In [21]:

```
1 #scatter plot
2 plt.scatter(df_train['GrLivArea'], df_train['SalePrice']);
```

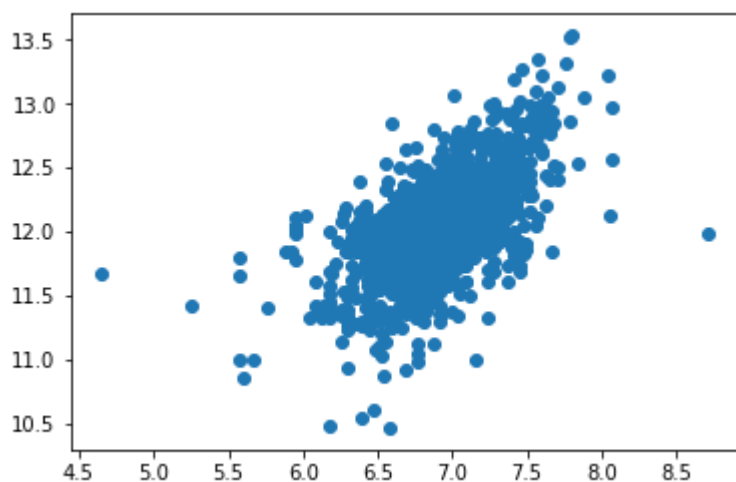


対数を取る前の古いバージョンの散布図は円錐型です（もどってみてみて！）今のグラフはそうじゃないだろ。これが正常性の力だ！ 正常性を担保すると均質性も解決してるのさ。

「販売価格」と「地下の敷地面積」を見てみよう

In [23]:

```
1 plt.scatter(df_train[df_train['TotalBsmtSF']>0]['TotalBsmtSF'], df_train[df_train['TotalBsmtSF']>
```



一般的に「販売価格」と「地下の敷地面積」は同じ分散レベルに収まっているだろう。クールだ！

最後にダミー変数

イージーモード

In [24]:

```
1 # カテゴリーデータをダミーにする
2 df_train = pd.get_dummies(df_train)
```

In [25]:

```
1 df_train.head()
```

Out[25]:

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

5 rows × 291 columns

結論

これでおしまい！最後まで来たよ！

このカーネルを通して、「Multivariate Data Analysis」の多くの戦略を練習してきたよ。変数とは何かを考え、「販売価格」だけに目を向けて分析して、関係する変数を分析して、欠損値やはずれ値も扱った。また本質的な統計的仮説も検証したし、カテゴリーデータをダミーにもしてきた。pythonが助けもあってこんなにたくさんの仕事をこなしてきたんだ。

でもまだ話は終わっていない。調査の直前でとまっている。「販売価格」さんに電話してディナーに招待しよう。そして彼女の振る舞いを予測しよう。彼女は正規化線形回帰アプローチを楽しむクチかい？それともアンサンブルメソッド、それとも他の手法？

それを見つけるのはあなた次第です。

データの傾向がわかったところで実際のモデル作成に取り掛かるうってこと？