

```
In [47]: from google.colab import files
uploaded = files.upload()

#import some necessary libraries

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", "."]).decode("utf8")) #check the files available in the directory

sample_data
submission.csv
test.csv
train.csv
```

```
In [0]: #Now Let's import and put the train and test datasets in pandas dataframe

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

```
In [0]: #display the first five rows of the train dataset.
train.head(5)
```

Out[0]:

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

5 rows × 11 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	Util
0	1461	20	RH	80.000	11622	Pave	NaN	Reg	Lvl	Al
1	1462	20	RL	81.000	14267	Pave	NaN	IR1	Lvl	Al
2	1463	60	RL	74.000	13830	Pave	NaN	IR1	Lvl	Al
3	1464	60	RL	78.000	9978	Pave	NaN	IR1	Lvl	Al
4	1465	120	RL	43.000	5005	Pave	NaN	IR1	HLS	Al

5 rows × 11 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' column 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))
```

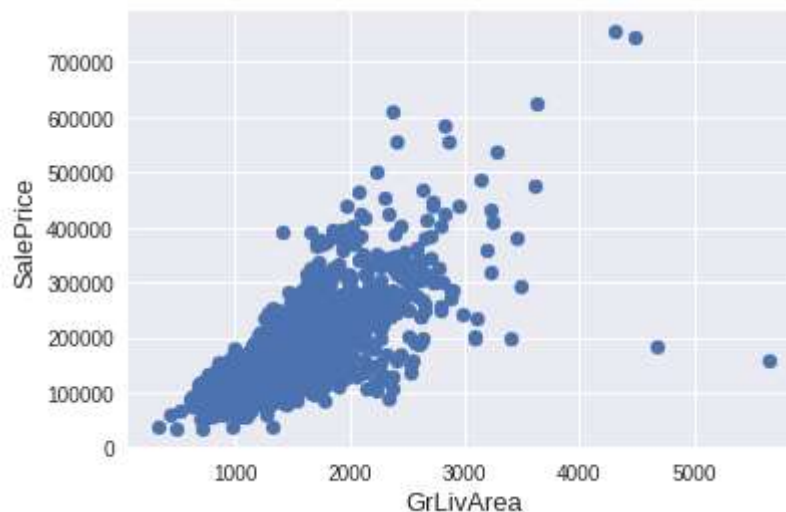
The train data size before dropping Id feature is : (1460, 81)

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

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

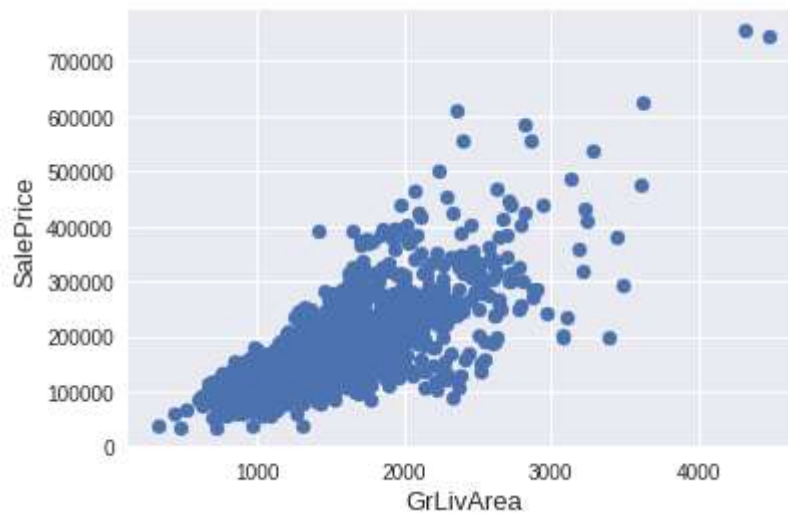
The test data size after dropping Id feature is : (1459, 79)

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



```
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()
```



```

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

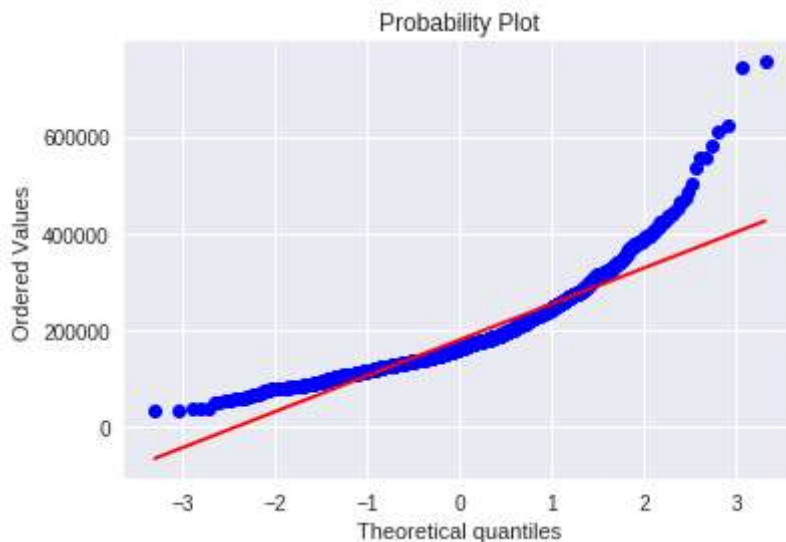
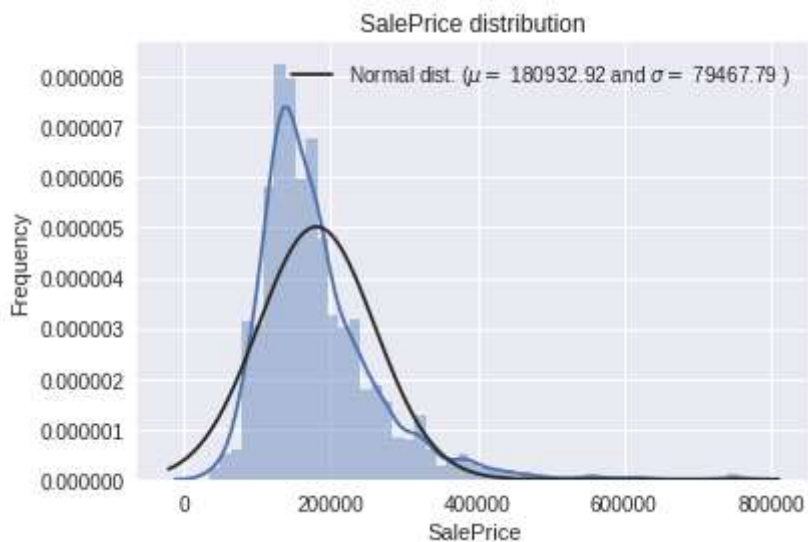
# Get the fitted parameters used by the function
(mu, sigma) = norm.fit(train['SalePrice'])
print( '\n mu = {:.2f} and sigma = {:.2f}\n'.format(mu, sigma))

#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



```
In [9]: #Log-transformation of the target variable

#We use the numpy fuction log1p which applies  $\log(1+x)$  to all elements of the
column
train["SalePrice"] = np.log1p(train["SalePrice"])

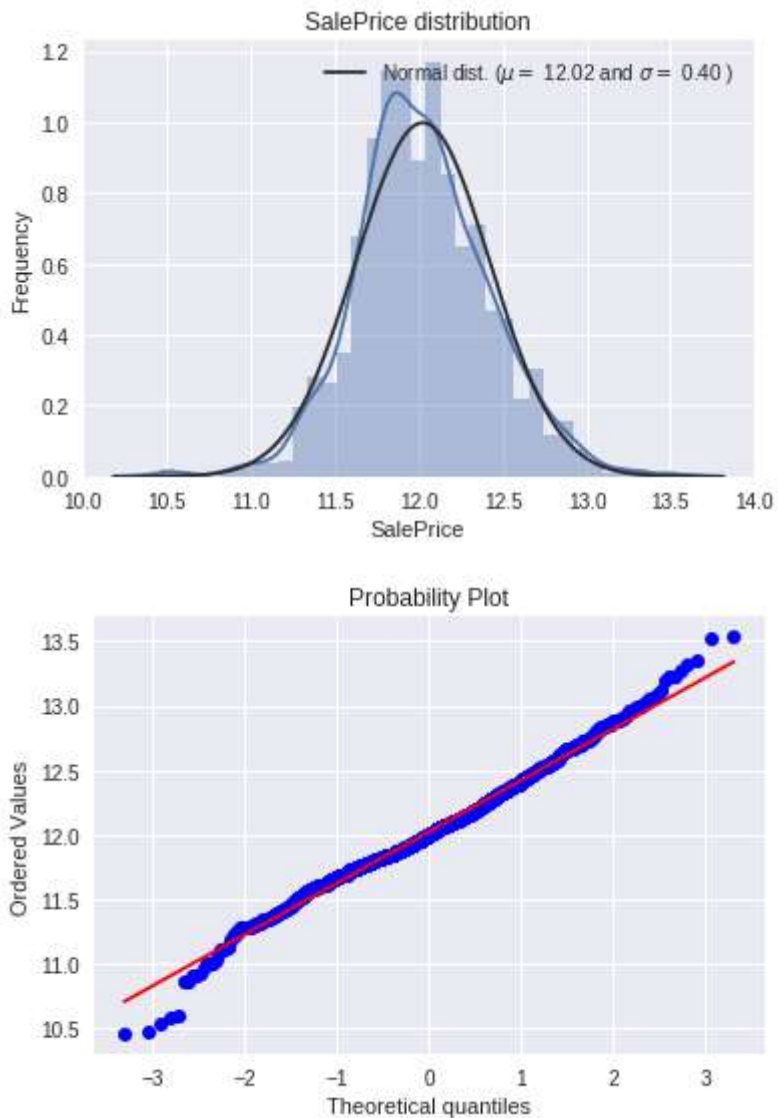
#Check the new distribution
sns.distplot(train['SalePrice'] , fit=norm);

# 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, si
gma)],
          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  $\sigma = 0.40$



```
In [10]: #Features engineering

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)
```

In [11]: *#Missing Data*

```
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_values(ascending=False)[:30]
missing_data = pd.DataFrame({'Missing Ratio' :all_data_na})
missing_data.head(20)
```

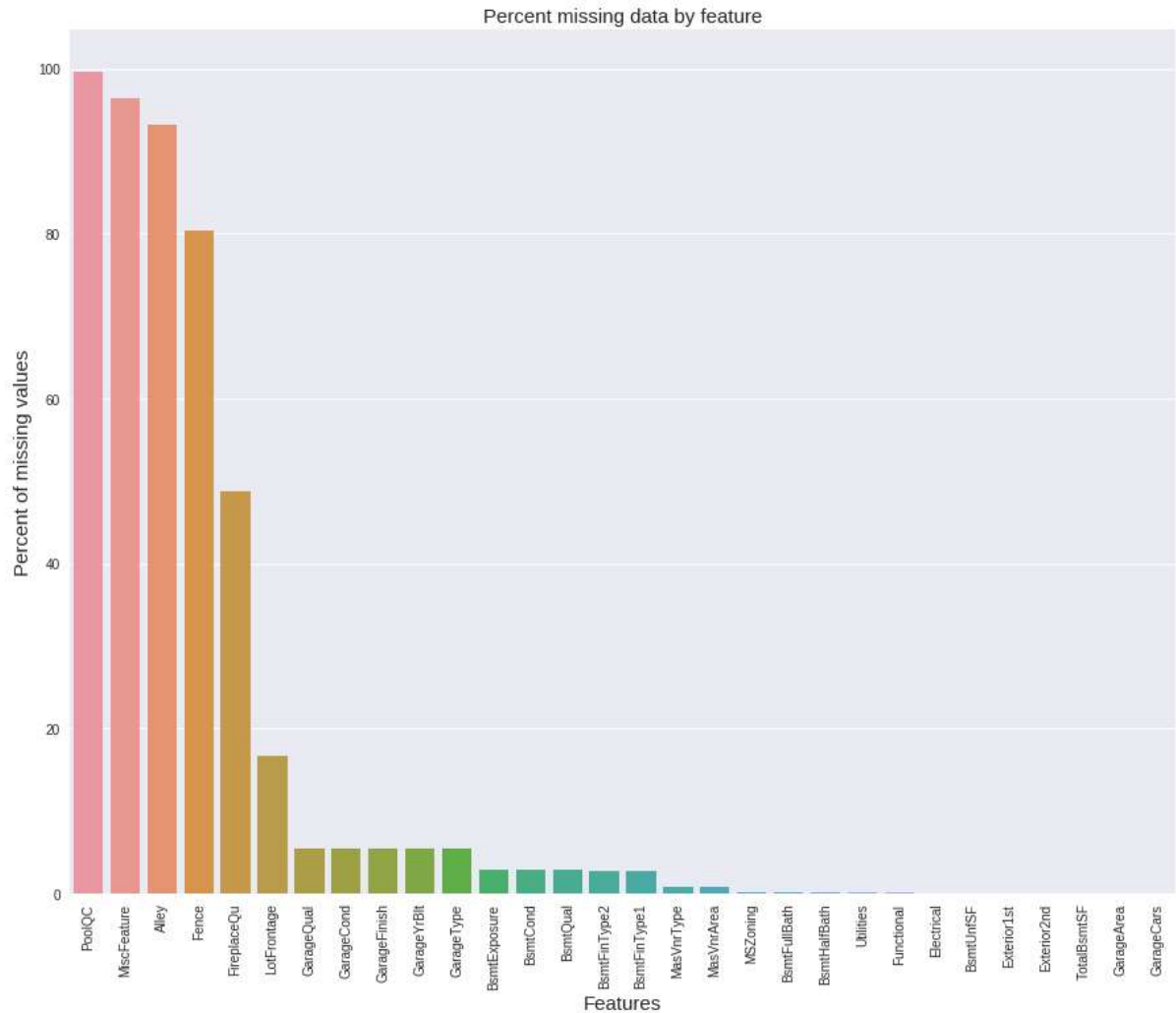
Out[11]:

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

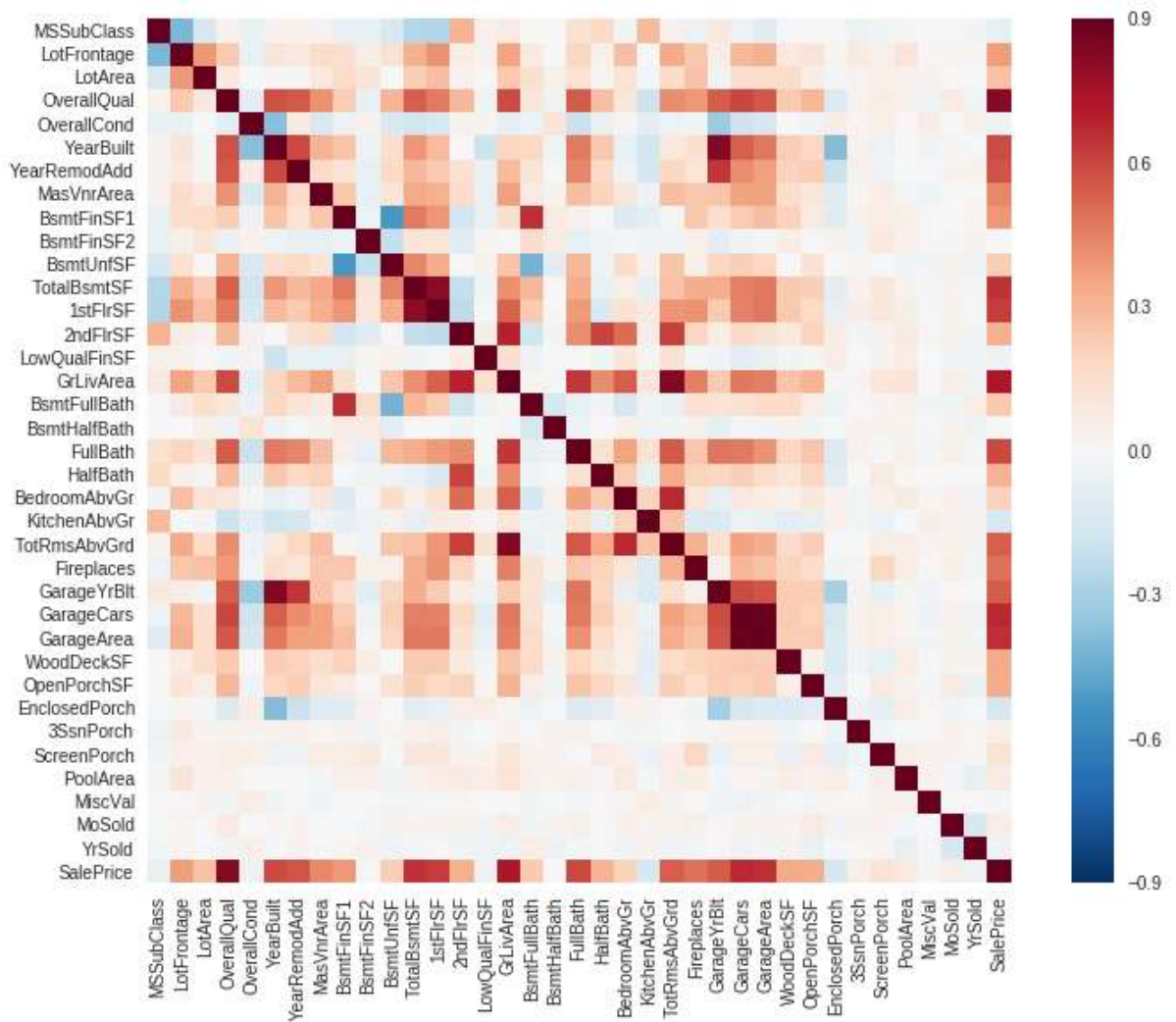
Out[12]: Text(0.5, 1.0, 'Percent missing data by feature')



In [13]: *#Data Correlation*

```
#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 0x7f819be34208>



```

In [0]: #Input missing values

all_data["PoolQC"] = all_data["PoolQC"].fillna("None")
all_data["MiscFeature"] = all_data["MiscFeature"].fillna("None")
all_data["Alley"] = all_data["Alley"].fillna("None")
all_data["Fence"] = all_data["Fence"].fillna("None")
all_data["FireplaceQu"] = all_data["FireplaceQu"].fillna("None")

#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()))

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

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

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

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

all_data["MasVnrType"] = all_data["MasVnrType"].fillna("None")
all_data["MasVnrArea"] = all_data["MasVnrArea"].fillna(0)

all_data['MSZoning'] = all_data['MSZoning'].fillna(all_data['MSZoning'].mode()[0])

all_data = all_data.drop(['Utilities'], axis=1)

all_data["Functional"] = all_data["Functional"].fillna("Typ")

all_data['Electrical'] = all_data['Electrical'].fillna(all_data['Electrical'].mode()[0])

all_data['KitchenQual'] = all_data['KitchenQual'].fillna(all_data['KitchenQual'].mode()[0])

all_data['Exterior1st'] = all_data['Exterior1st'].fillna(all_data['Exterior1st'].mode()[0])
all_data['Exterior2nd'] = all_data['Exterior2nd'].fillna(all_data['Exterior2nd'].mode()[0])

all_data['SaleType'] = all_data['SaleType'].fillna(all_data['SaleType'].mode()[0])

all_data['MSSubClass'] = all_data['MSSubClass'].fillna("None")

```

```
In [15]: #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_values(ascending=False)
missing_data = pd.DataFrame({'Missing Ratio' :all_data_na})
missing_data.head()
```

Out[15]:

Missing Ratio

```
In [0]: #More features engeneering

#Transforming some numerical variables that are really categorical

#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)
```

```
In [17]: #Label Encoding some categorical variables that may contain information in the
ir ordering set

from sklearn.preprocessing import LabelEncoder
cols = ('FireplaceQu', 'BsmtQual', 'BsmtCond', 'GarageQual', 'GarageCond',
        'ExterQual', 'ExterCond', 'HeatingQC', 'PoolQC', 'KitchenQual', 'BsmtFin
nType1',
        'BsmtFinType2', 'Functional', 'Fence', 'BsmtExposure', 'GarageFinish',
        'LandSlope',
        'LotShape', 'PavedDrive', 'Street', 'Alley', 'CentralAir', 'MSSubClas
s', '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)

```
In [0]: # Adding total sqfootage feature
all_data['TotalSF'] = all_data['TotalBsmtSF'] + all_data['1stFlrSF'] + all_dat
a['2ndFlrSF']
```

In [19]: *#Skewed features*

```
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[19]:

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

In [20]: *#Box Cox Transformation of (highly) skewed features*

```
skewness = skewness[abs(skewness) > 0.75]
print("There are {} skewed numerical features to Box Cox transform".format(ske
wness.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

In [21]: *#Getting dummy categorical features*

```
all_data = pd.get_dummies(all_data)
print(all_data.shape)

(2917, 220)
```

In [0]: *#Getting the new train and test sets*

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

## Modelling

In [0]: *#Import librairies*

```
from sklearn.linear_model import ElasticNet, Lasso, BayesianRidge, LassoLarsI
C
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
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
```

In [0]: *#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="neg_
mean_squared_error", cv = kf))
    return(rmse)
```

```
In [0]: lasso = make_pipeline(RobustScaler(), Lasso(alpha =0.0005, random_state=1))
ENet = make_pipeline(RobustScaler(), ElasticNet(alpha=0.0005, l1_ratio=.9, random_state=3))
KRR = KernelRidge(alpha=0.6, kernel='polynomial', degree=2, coef0=2.5)
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)
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)
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 = 1
1)
```

```
In [26]: score = rmsle_cv(lasso)
print("\nLasso score: {:.4f} ({:.4f})\n".format(score.mean(), score.std()))

Lasso score: 0.1115 (0.0074)
```

```
In [27]: score = rmsle_cv(ENet)
print("ElasticNet score: {:.4f} ({:.4f})\n".format(score.mean(), score.std()))

ElasticNet score: 0.1116 (0.0074)
```

```
In [28]: score = rmsle_cv(KRR)
print("Kernel Ridge score: {:.4f} ({:.4f})\n".format(score.mean(), score.std()))

Kernel Ridge score: 0.1153 (0.0075)
```

```
In [29]: score = rmsle_cv(GBoost)
print("Gradient Boosting score: {:.4f} ({:.4f})\n".format(score.mean(), score.std()))

Gradient Boosting score: 0.1177 (0.0080)
```

```
In [30]: score = rmsle_cv(model_xgb)
print("Xgboost score: {:.4f} ({:.4f})\n".format(score.mean(), score.std()))

Xgboost score: 0.1164 (0.0069)
```

```
In [31]: score = rmsle_cv(model_lgb)
print("LGBM score: {:.4f} ({:.4f})\n".format(score.mean(), score.std()))
```

LGBM score: 0.1167 (0.0072)

```
In [0]: # Stacking models

# Averaging base models

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)
```

```
In [33]: # Averaged base models score

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)



```

In [0]: # Stacking averaged Models Class

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_models_)))

        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
        # a 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 averaged
    # 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(
axis=1)
                for base_models in self.base_models_ ])
        return self.meta_model_.predict(meta_features)

```

```

In [37]: # Stacking Averaged models Score

stacked_averaged_models = StackingAveragedModels(base_models = (ENet, GBoost,
KRR),
                                                    meta_model = lasso)

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

```

Stacking Averaged models score: 0.1084 (0.0074)

```
In [0]: def rmsle(y, y_pred):  
        return np.sqrt(mean_squared_error(y, y_pred))
```

## Final Training and Prediction

```
In [40]: # StackedRegressor  
stacked_averaged_models.fit(train.values, y_train)  
stacked_train_pred = stacked_averaged_models.predict(train.values)  
stacked_pred = np.expm1(stacked_averaged_models.predict(test.values))  
print(rmsle(y_train, stacked_train_pred))  
  
0.07814905151549066
```

```
In [41]: # XGBoost  
model_xgb.fit(train, y_train)  
xgb_train_pred = model_xgb.predict(train)  
xgb_pred = np.expm1(model_xgb.predict(test))  
print(rmsle(y_train, xgb_train_pred))  
  
0.07898668813465673
```

```
In [42]: # LightGBM  
model_lgb.fit(train, y_train)  
lgb_train_pred = model_lgb.predict(train)  
lgb_pred = np.expm1(model_lgb.predict(test.values))  
print(rmsle(y_train, lgb_train_pred))  
  
0.07343743130986105
```

```
In [43]: '''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:  
0.07548543390046325
```

```
In [0]: # Ensemble prediction  
ensemble = stacked_pred*0.70 + xgb_pred*0.15 + lgb_pred*0.15
```

```
In [0]: # Submission  
sub = pd.DataFrame()  
sub['Id'] = test_ID  
sub['SalePrice'] = ensemble  
sub.to_csv('submission.csv', index=False)
```

```
In [0]: # Download submission file  
files.download('submission.csv')
```