# Project 2
# Kaggle Competition
# House Prices: Advanced Regression Techniques

***Team Members***
**Karthik Karunanithi (802827527)**
**Bharath Krishnan (893429449)**
**Shankar Tiwari (803012350)**
**Gargi Mrunal Kulkarni (893210922)**

# TABLE OF CONTENTS

# 1 INTRODUCTION

In this project we signed up for the Kaggle competition. The aim was to complete the competition and submit the code for rank. Then improve the code for better analysis to improve the rank. The competition we entered was "House Prices: Advanced Regression Techniques". The competition details can be found on site: https://www.kaggle.com/c/house-prices-advanced-regression-techniques. We implemented different advanced regression algorithms and compared the results for best prediction.

Following are the details for the project implementations:

**Dataset**: Provided by Kaggle and in known as Ames Housing Dataset

**Data Mining Tool**: Python scikit library.

**Analysis & Prediction:**

Prediction of the sale price of the houses

**Algorithms**: The following algorithms were implemented in the project:

Advanced Regression Techniques like LASSO, XgBoost, PCA etc.

# 2 ABOUT THE DATASET AND ALGORITHMS

## 2.1 DATASET

The dataset was provided by Kaggle. It is Ames Housing Dataset. It's an alternative, modernized and expanded version, of the often cited Boston Housing dataset. The details are available on site https://www.kaggle.com/c/house-prices-advanced-regression-techniques/data.

## 2.2 ALGORITHMS

The algorithms implemented to predict sale price using Ames Housing Dataset are as given below :

1. Lasso Regression: In statistics and machine learning, lasso (least absolute shrinkage and selection operator) (also Lasso or LASSO) is a regression analysis method that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretability of the statistical model it produces. Lasso was originally formulated for least squares models and this simple case reveals a substantial amount about the behavior of the estimator, including its relationship to ridge regression and best subset selection and the connections between lasso coefficient estimates and so-called soft thresholding. It also reveals that (like standard linear regression) the coefficient estimates need not be unique if covariates are collinear. Though originally defined for least squares, lasso regularization is easily extended to a wide variety of statistical models including generalized linear models, generalized estimating equations, proportional hazards models, and M-estimators, in a straightforward fashion. Lasso's

ability to perform subset selection relies on the form of the constraint and has a variety of interpretations including in terms of geometry, Bayesian statistics, and convex analysis.

2. Xgboost: Xgboost is an open-source software library which provides the Gradient boosting framework for C++, Java, Python, R, and Julia. It works on Linux, Windows, and macOS. From the project description, it aims to provide a "Scalable, Portable and Distributed Gradient Boosting (GBM, GBRT, GBDT) Library". Other than running on a single machine, it also supports the distributed processing frameworks Apache Hadoop, Apache Spark, and Apache Flink. It has gained much popularity and attention recently as it was the algorithm of choice for many winning teams of a number of machine learning competitions like Kaggle.

3. Ridge Regression: Similar to LASSO, ridge regression is a regression analysis method that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretability of the statistical model it produces. Linear regression models are often fitted using the least squares approach, but they may also be fitted in other ways, such as by minimizing the "lack of fit" in some other norm (as with least absolute deviations regression), or by minimizing a penalized version of the least squares loss function as in ridge regression (L2-norm penalty) .

4. Random Forest regression: Random forests or random decision forests are an ensemble learning method for classification, regression and other tasks, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random decision forests correct for decision trees' habit of overfitting to their training set.

5. Principal Component Analysis(PCA): Principal component analysis (PCA) is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components. The number of principal components is less than or equal to the number of original variables. This transformation is defined in such a way that the first principal component has the largest possible variance (that is, accounts for as much of the variability in the data as possible), and each succeeding component in turn has the highest variance possible under the constraint that it is orthogonal to the preceding components. The resulting vectors are an uncorrelated orthogonal basis set. PCA is sensitive to the relative scaling of the original variables. PCA is mostly used as a tool in exploratory data analysis and for making predictive models. PCA can be done by eigenvalue decomposition of a data covariance (or correlation) matrix or singular value decomposition of a data matrix, usually after mean centering (and normalizing or using Z-scores) the data matrix for each attribute.[4]The results of a PCA are usually discussed in terms of component scores, sometimes called factor scores (the transformed variable values corresponding to a particular data point), and loadings (the weight by which each standardized original variable should be multiplied to get the component score). In regression analysis, the larger the number of explanatory variables allowed, the greater is the chance of overfitting the model, producing conclusions that fail to generalise to other datasets. One approach, especially when there are strong correlations between different possible explanatory variables, is to reduce them to a few principal components and then run the regression against them, a method called principal component regression.

# 3 ABOUT THE TOOL – PYTHON SKLEARN

Scikit-learn (formerly scikits.learn) is a free software machine learning library for the Python programming language. It features various classification, regression and clustering algorithms including support vector machines, random forests, gradient boosting, k-means and DBSCAN, and is designed to interoperate with the Python numerical and scientific libraries NumPy and SciPy.

scikit-learn is known as Machine Learning in Python with following features:
Simple and efficient tools for data mining and data
analysis Accessible to everybody, and reusable in
various contexts Built on NumPy, SciPy, and Matplotlib
Open source, commercially usable - BSD license

## 3.1 SYSTEM REQUIREMENTS AND INSTALLATION STEPS

The system requirement for Python SKLearn have no any minimal specification. Since data analysis is a computationally intensive task—the better your hardware, the better your experience. Also, the memory should be enough to handle big data sets.

The installation steps of SKLearn is given on the site: [http://scikit-learn.org/stable/install.html](http://scikit-learn.org/stable/install.html) Other installation required to support the scikit learn library are:
Python (>= 2.6 or >= 3.3) - Python version of project is
version 2.6 NumPy (>= 1.6.1)
SciPy (>= 0.9)
Note: The Python version required is 2.6. If version is different and libraries installed are for different version, code will fail to execute. For correct installation please check online instructions.

The steps to run the code are as follows:
Install Python 2.6 and sklearn library, sciPy, numPy, matplotlib and
pandas Download the dataset form [here](here) and save it on local drive.
Copy the code file on the local drive.
Open the file and change the dataset path (Test and Train) to the data files saved on the local machine.

```python
# Reading the datasets

train_df =
pd.read_csv('/home/bharathkrishnan/bharath/Pro
2/train.csv', index_col=0)
test_df =
pd.read_csv('/home/bharathkrishnan/bharath/Pro
2/test.csv', index_col=0)
```

Open command window and change the directory to one where code file is located. Run the code to produce output CSV file as required by the competition.
The algorithms in the code have been commented out. In order to run a particular algorithm, remove the comments from the respective code and assign the respective output variable to the variable y_final.

# 4 PROJECT IMPLEMENTATION

In this project, we implemented different regression techniques on Ames Housing Dataset and picked the one giving best predictions.

## 4.1 DATA PRE-PROCESSING

Most of the real world data are generally Incomplete, noisy or inconsistent. As a result, a certain set of procedures are followed to make the data fit for the analysis. The following measures are considered in the project:

1. Categorical Variables are coded using the dummy variables. The number of dummy variables for a categorical feature is equal to (Number of categories-1). Pandas, a Python Library, allows you to code the categorical variables into the dummy variables using the function,

**Dataframe.*get_dummies()***

2. Handling the missing Variables: Missing variables can be handled in different ways: Substituting it with the attribute mean, ignore the observation corresponding to the the missing value, Use the Label of the missing value as the target variable. Since the number of rows are significantly low, removing any number of observation will lead to a significant amount of information in the data. As a result, the missing values are substituted with the attribute mean using the following function.
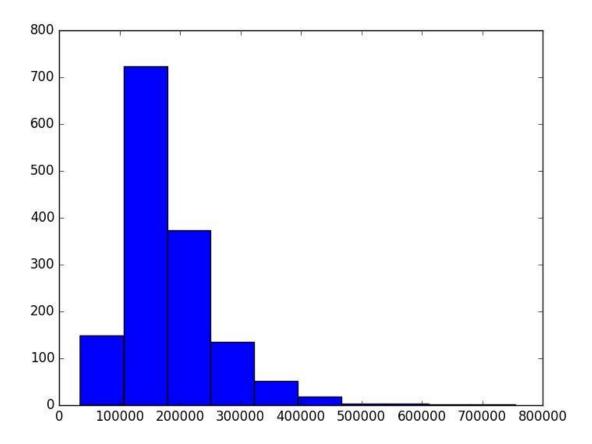
**Dataframe.*isnull*(Dataframe.*mean*())**

3. Handling the skewness: The data provided was skewed to a large extent. Analysis on a skewed data will lead to incorrect results or inference. The most commonly method used to reduce the skewness of the data is the *log* transformation.

4. Data Normalization: Data is normalized as the features in the data are generally of different unit. the intention is that these **normalized values** allow the comparison of corresponding normalized values for different datasets in a way that eliminates the effects of certain gross influences. Normalization is done using the following formula:
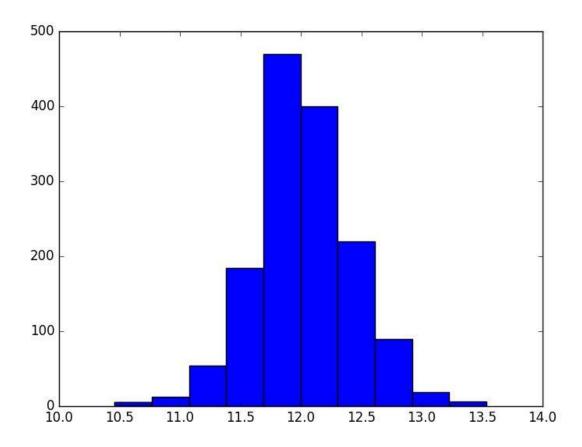
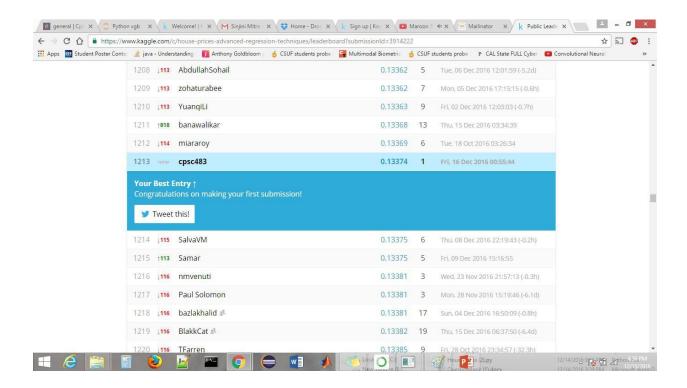$$\frac{X - \mu}{\sigma}$$

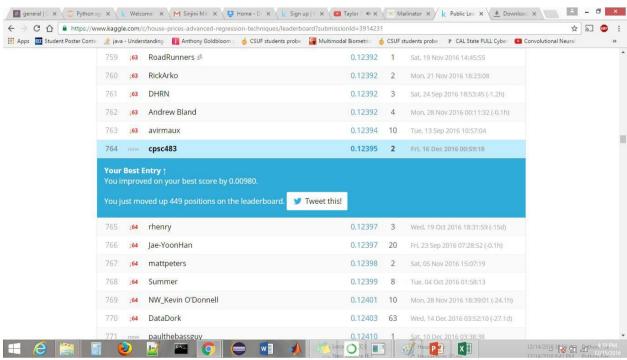## 4.2 SCREENSHOTS

1. Distribution of the data
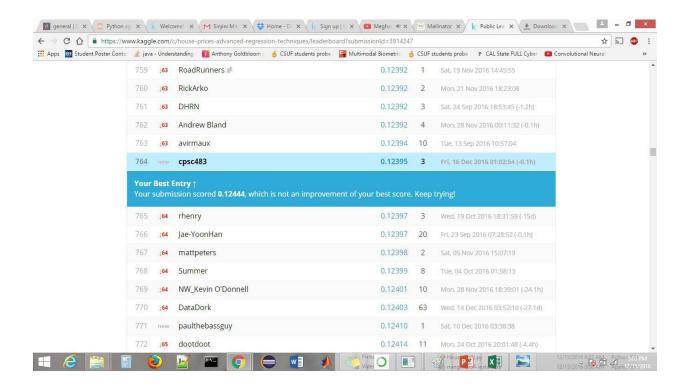
2. Data after Log Transformation

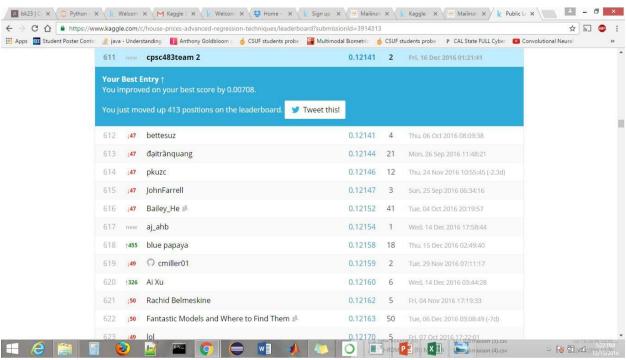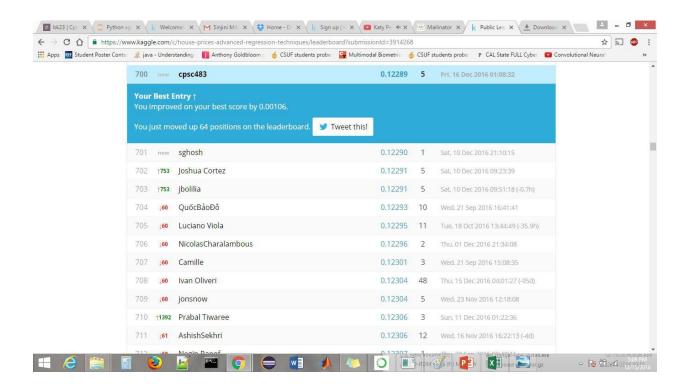3. Result of performing Principal Component Analysis and Lasso

## 4. Lasso



## 5. Ridge

## 6. Ridge Lasso and Xgboost



## 7. XGBoost and Lasso

# 5 REFERENCES

[1] http://scikit-learn.org/stable/index.html

[2] https://www.wikipedia.org/

[3] https://www.kaggle.com/c/house-prices-advanced-regression-techniques

[4] https://www.kaggle.com/apapiu/house-prices-advanced-regression-techniques/regularized-linear-models

[5] https://www.kaggle.com/humananalog/house-prices-advanced-regression-techniques/xgboost-lasso

[6] https://en.wikipedia.org/wiki/Normalization_(statistics)