# STATS 790
# Written Report for Final Project (Draft)

Amandeep Sidhu (400076920)

09 April, 2023

## Contents

## Introduction

In this final project, the data set used for our techniques will be the built in "Boston Housing" data set that is available in R and Python (Harrison Jr. and Rubinfeld 1978). According to the description provided by the original data collector, it was collected in the 1970s by Harrison and Rubinfeld during the census at this time. This data set includes 506 observations with 13 possible predictors that aim to predict the value of the houses in the area of Boston. Some of these predictors include: CRIM (crime rate), AGE (proportion of owner-occupied units built prior to 1940), DIS (weighted distances to employment centres) and LSTAT (the lower status/class of the population in Boston). The target variable (denoted as 'MEDV') is the median value of the homes in thousands of dollars. This data set was first used in a paper titled "Housing Prices and the Demand for Clean Air" (Harrison Jr. and Rubinfeld 1978) as the demand for clean air rose with pollution becoming a predominant issue in modern society. As a naive approach, Harrison and Rubinfeld used a simple linear regression model (least squares method) to estimate the parameters. On top of this, they simply used a variety of different combinations of the predictors present in attempt to improve model performance. The final model they used included all 13 predictor variables which yielded an adjusted $R^2$ value of 0.728. Remember, that the adjusted $R^2$ essentially is the amount of variation in the target variable explained by the model.

The statistical problem at hand here is to use regression to predict the house prices for houses in Boston. More specifically, we will implement Ridge Regression, LASSO and Elastic Net techniques on this dataset to address this regression problem in both R and Python to try and compare their time and memory for a series of benchmark problems. In other words, we will benchmark each algorithm in R and Python to compare how fast each technique is and the relative accuracy (using RMSE, $R^2$ and MAE) to discuss potential trade-offs between the techniques. It is important to note that although a technique may be less computationally expensive, it may not be worth a decline in accuracy (the opposite holds true as well).

As opposed to traditional data science courses, we will also be extending this application by trying to implement a warm start for those algorithms discussed above in order to obtain the optimal parameters for each model.

## Methodology

Before we dive into the specifics of our three model approaches, we must briefly discuss what is meant by a "warm start" approach. As discussed by (Chu et al., n.d.), we observe that a warm start is essentially a technique that aims to reduce the running time of the iterative methods in its model fitting process. It uses the solution of a different optimization as a starting point (rather than a cold approach that has no defined starting point). Ideally, one would choose a starting point that is very close to the optimal point, but generally multiple starting points are investigated.

### Method 1: Ridge Regression

The first method to be implemented on this data set is called Ridge Regression. Ridge regression is a model tuning method that can be used to analyze data and performs L2 regularization (G. L. Team 2022). The estimates equation for ridge regression is then: $\hat{\beta^{ridge}} = argmin_\beta\{\sum_{i=1}^{N}(y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j)^2 + \lambda\sum_{j=1}^{p}|\beta_j|^2\}$ (Hastie, Tibshirani, and Friedman 2008). The penalty term aims to prevent multicollinearity, reduce model complexity by coefficient shrinkage. As a result, we will need to tune this specific parameter to find the optimal value as this will significant impact the results of our model fitting. In order to use a warm start approach, the ridge function has a 'start' argument that aims to do this by using a previously defined ridge regression model. This process will be repeated in Python as well.

### Method 2: LASSO

The second method to be implemented on this data set is called LASSO regression. Like ridge regression, LASSO regression is also a regularization technique. The primary point of this technique is that it uses shrinkage which corresponds to shrinking the data values towards a central point as a mean (Kumar 2023). As this article also points out, the LASSO procedure encourages simpler models (i.e. with fewer parameters present). However, as opposed to ridge regression, LASSO uses L1 regularization techniques. The estimates equation for LASSO is then: $\hat{\beta^{lasso}} = argmin_\beta\{\frac{1}{2}\sum_{i=1}^{N}(y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j)^2 + \lambda\sum_{j=1}^{p}|\beta_j|\}$ (Hastie, Tibshirani, and Friedman 2008). The term $\lambda$ here indicates the amount of shrinkage present. This will be tuned to attain optimal estimates. Once again, within the 'glmnet' package, the lasso function can be defined to have a 'start' argument of a previously defined model. This will be the starting point for our warm

start approach. This process will be repeated in Python as well.

**Method 3: Elastic Net**

The final method to be implemented on this data set is called Elastic Net. In general terms, Elastic Net regression uses penalties from both the LASSO and ridge techniques to try and regularize regression models (C. Team 2022). In other words, it tries to learn from both LASSO and ridge regression to try and improve on its own model creation. This technique may seem rather redundant to do given that both LASSO and ridge regression are already being implemented, however we wish to also observe the timing and memory storage differences across multiple programming languages. A method that combines both LASSO and Ridge regression may be pose interesting results for this reason. Moreover, as C. Team (2022) explains, elastic net improves on LASSO's limitations through the inclusion of "n" number of variables until saturation. Once again, a warm start will be implemented to try and find the optimal point. This process will be repeated in Python as well.

**Description of Comparison Criterion**

As this is a regression problem by nature, standard criterion such as RMSE, $R^2$ and even MAE will be calculated to compare the three methods. In addition, we will also use timing metrics (such as R's benchmark tool) in order to find the computation time for the algorithms to run. This will also be computed in Python as well.

**Results**

To avoid clutter in this report, I have included the coding portion in the appendix of the report.

**Exploratory Data Analysis (EDA)**

Before we actually fit the model using the methods discussed, let us first explore the data in order to see if any pre-processing or scaling is required. Upon initial inspection, our target variable is 'MEDV' and the remaining 13 variables will serve as the predictors in model creation. To observe the first six rows of the data set, refer to the Figure **1** in the appendix section of the report. Overall, the values appear to be normal for the context of their variable (i.e. the nitrogen counts are appropriate and not randomly in the millions). There appears to also be no missing values as well. To observe the spread of the variables, we will observe the boxplot.
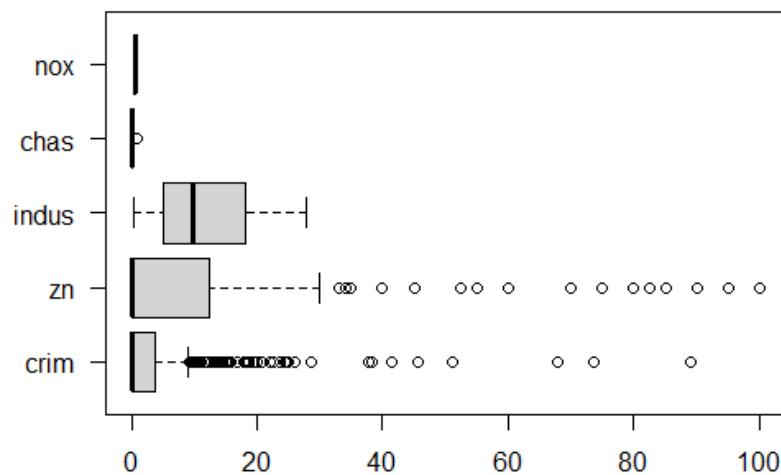


**Figure 2:** Boxplot of the first five variables present

For the most part, there is a good spread for the variables as a whole as the box manages to capture them adequately. We have only included the first five variables as the boxplot would be too crowded to interpret (in the appendix, simply change the code to observe the other variables spread for the

5

boxplot). Overall, there may be an issue in terms of outliers for the crim and zn predictors, but we will keep them in for now as we are unclear on how they may affect other predictors.

In order to observe the effect of the predictors on each other, let us create a correlation plot to observe potential relationships between the variables. Figure **3** shows a correlation plot between the variables present.
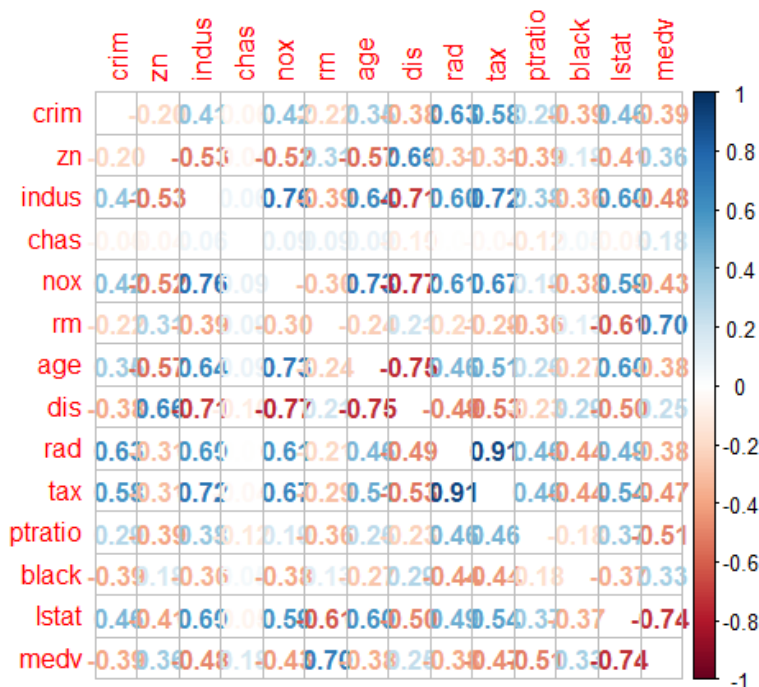


**Figure 3:** Correlation Matrix of the Variables Present

Obsevre that a darker shade of blue indicates a stronger positive relationship, while a darker red indicates a stronger negative relationship. Some observations of this plot include, the crime variable is strongly correlated with the variables rad and tax. This suggests that if either of these variables increase then so too will the crime variable. Moreover, if we observe the nox variable, we see this is strongly associated with the indus in a positive relationship but a negative relationship with the distance (of employment centres). Remember that the "indus" variable denotes the industrial area so it would make sense to have a higher nitrogen oxide count there. As discussed before, we will keep the variables with potential outliers in the data set as they seem to have a strong correlation with other variables in the data set.

Now, let us visualize the data with the help of a scatterplot of dependent variables versus the

median value of houses (i.e. our "MEDV" attribute). As discussed by Subrahmanya (2018), we can use 'ggplot2' to attain this important visual.
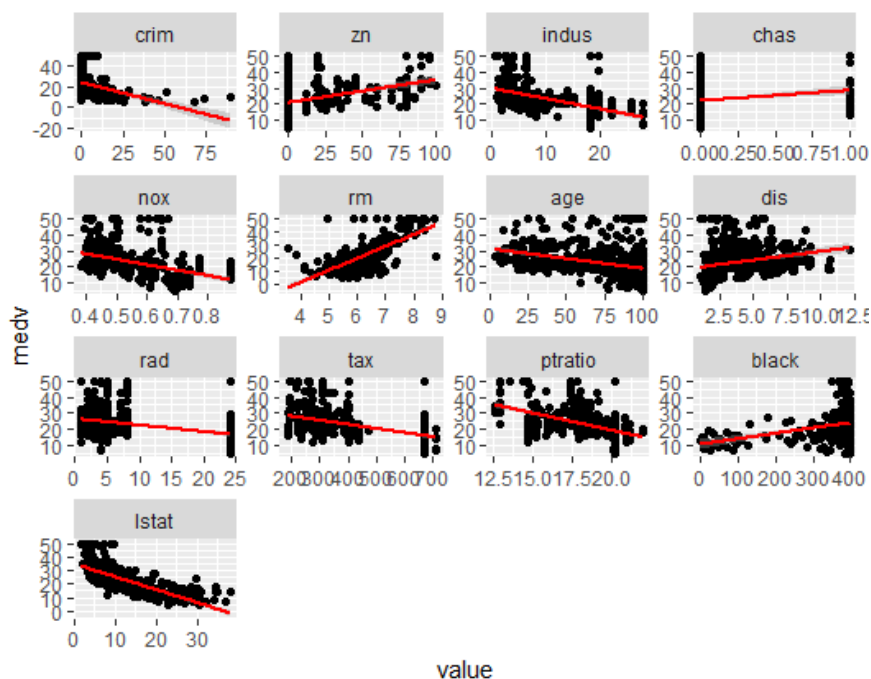


**Figure 4:** Scatterplot versus the MEDV present

As we can see in Figure **4**, we can observe the relationship between a single predictor variable and the MEDV variable (i.e. the target variable). Some observations include that the per capita crime (crim) has a negative relationship with the MEDV target variable, while the rm variable has a positive relationship with MEDV. These results are also in agreement of our correlation matrix in Figure **3**, but that figure gives the numerical value of the slope.

**Results of the Regression Techniques in R and Python**

Once again, to observe the coding portion of this project, refer to the appendix section for this.

| Model Implemented | RMSE | MAE | $R^2$ | Time (seconds) |
| --- | --- | --- | --- | --- |
| Ridge (in R) | 4.674069 | 3.435065 | 0.7317628 | 0.05 |
| Lasso (in R) | 4.71244 | 3.46451 | 0.7276256 | 0.07 |
| Elastic Net (in R) | 4.723965 | 3.474352 | 0.7264075 | 0.06 |
| Ridge (in Python) | 4.94192 | 3.504674 | 0.694377 | 0.29501 |
| Lasso (in Python) | 4.941965 | 3.504644 | 0.6943718 | 0.2730556 |
| Elastic Net (in Python) | 4.943073 | 3.504854 | 0.69423475 | 0.26695085 |

Table 1: Results of Model Implementation

## Conclusions

In terms of model performance (i.e. observing the metrics in Table **1**), it appears that all of them do relatively similar. In both R and Python, we observe that the RMSE settles around 4.7 (which corresponds to the standard deviation of the residuals) and because this is low we know that the model is fit adequately. We observe that the MAE settles around 3.5 and again a smaller value here is good. Next, $R^2$ seems to have an overall score of around 0.7 which explains variation in the data by the model. The subtle differences between R and Python may be explained when we split the data. Though we set the seed the same, the data may not all be equal in terms of the train and test split.

As we can again see in Table **1**, the metrics appear to be relatively close together (with small differences moving from Python to R). This may be the result of a different parameter definition, as we tried to keep the optimization process the same across both programming languages. However, observe the difference in the time it takes to compute the process. The times in Python appear to be longer than those in R and this may be the result of the built in 'warm_start' for R versus explicitly defining the process in Python. Ridge Regression and Elastic Net do not have a warm_start argument and thus we had to explicitly define this process. As a result, we observe a longer time. Thus, we conclude that although we are testing the same data present, we cannot expect to observe exact results. That said, it is essentially the difference of less than a second so an argument can be made that this is insignificant. Overall, the timing of these algorithms are not the same when implementing a warm_start approach.

# References

Chu, Bo-Yu, Chia-Hua Ho, Cheng-Hao Tsai, Chieh-Yen Lin, and Chih-Jen Lin. n.d. *Warm Start for Parameter Selection of Linear Classifiers.* https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/warm-start/warm-start.pdf.

Harrison Jr., David, and David L. Rubinfeld. 1978. *Hedonic Housing Prices and the Demand for Clean Air.* https://www.law.berkeley.edu/files/Hedonic.PDF.

Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2008. *The Elements of Statistical Learning.* Springer.

Kumar, Dinesh. 2023. *A Complete Understanding of LASSO Regression.* https://www.mygreatlearning.com/blog/understanding-of-lasso-regression/.

Subrahmanya, Rashmi. 2018. *Analysis of Boston Housing Data Using Linear Regression, Trees and GAM.* https://rpubs.com/Rashmi_Subrahmanya/371719.

Team, CFI. 2022. *Elastic Net- a Regression Method That Performs Variable Selection and Regularization Simultaneously.* https://corporatefinanceinstitute.com/resources/data-science/elastic-net/.

Team, Great Learning. 2022. *What Is Ridge Regression?* https://www.mygreatlearning.com/blog/what-is-ridge-regression/.

# Appendix

```
library(MASS)
library(corrplot)
library(dplyr)
library(ggplot2)
library(glmnet)
library(caret)
data("Boston")
head(Boston)
```

```
##      crim zn indus chas   nox    rm  age    dis rad tax ptratio  black lstat
## 1 0.00632 18  2.31    0 0.538 6.575 65.2 4.0900   1 296    15.3 396.90  4.98
## 2 0.02731  0  7.07    0 0.469 6.421 78.9 4.9671   2 242    17.8 396.90  9.14
## 3 0.02729  0  7.07    0 0.469 7.185 61.1 4.9671   2 242    17.8 392.83  4.03
## 4 0.03237  0  2.18    0 0.458 6.998 45.8 6.0622   3 222    18.7 394.63  2.94
## 5 0.06905  0  2.18    0 0.458 7.147 54.2 6.0622   3 222    18.7 396.90  5.33
## 6 0.02985  0  2.18    0 0.458 6.430 58.7 6.0622   3 222    18.7 394.12  5.21
##   medv
## 1 24.0
## 2 21.6
## 3 34.7
## 4 33.4
## 5 36.2
## 6 28.7
```

**Figure 1:** First six rows of the Boston Housing Data Set

```
#code for the correlation matrix
corrplot(cor(Boston), method = "number", diag = FALSE)
```

```r
#code for creating the boxplot of the first five variables
boxplot(Boston[c(1,2,3,4,5)], horizontal= TRUE, las= 1)
```

```r
#code for the scatterplots vs dependent variable


#reshape the data
data <- reshape2::melt(Boston, id.vars = "medv")


#create the plot
ggplot(data, aes(x = value, y = medv)) +
  geom_point() +
  stat_smooth(method = "lm", se = TRUE, col = "red") +
  facet_wrap(~ variable, scales = "free") +
  theme_gray()
```

```r
#code for implementation of ridge regression


set.seed(400076920)


#Split into 75-25 accordingly
train_ind <- createDataPartition(Boston$medv, p = 0.75, list = FALSE)
train <- Boston[train_ind,]
test <- Boston[-train_ind,]


#define predictor and response/target variable for the cv.glment
 x<- model.matrix(medv ~. , data = train)[,-1]
 y <- train$medv


 #create the ridge regression with a warm start/ and find timing
 time_ridge <-system.time({
   ridge_model <- cv.glmnet(x,y, alpha=0, nfolds = 10,
                            lambda=seq(0.001, 0.1, length=100), standardize=TRUE,
```

```r
                        keep=TRUE, type.measure="mse", warm.start=TRUE)


})


#now, use the model to make predictions
x_test <- model.matrix(medv~., data = test)[,-1]
y_test <- test$medv
y_pred <- predict(ridge_model, newx = x_test, s = ridge_model$lambda.min)


#calculate comparison criterion


rmse_ridge <- sqrt(mean((y_pred-y_test)^2))
mae_ridge <- mean(abs(y_pred - y_test))
rsq_ridge <- cor(y_pred, y_test)^2


#code implementation for lasso regression
set.seed(400076920)


#Split into 75-25 accordingly
train_ind <- createDataPartition(Boston$medv, p = 0.75, list = FALSE)
train <- Boston[train_ind,]
test <- Boston[-train_ind,]


#define predictor and response/target variable for the cv.glment
 x<- model.matrix(medv ~. , data = train)[,-1]
 y <- train$medv


 #create the lasso regression with a warm start/ and find timing
 time_lasso <-system.time({
   lasso_model <- cv.glmnet(x,y, alpha=1, nfolds = 10,
                        lambda=seq(0.001, 0.1, length=100), standardize=TRUE,
                        keep=TRUE, type.measure="mse", warm.start=TRUE)
```

```r
})


 #now, use the model to make predictions
x_test <- model.matrix(medv~., data = test)[,-1]
y_test <- test$medv
y_pred <- predict(lasso_model, newx = x_test, s = lasso_model$lambda.min)


#calculate comparison criterion


rmse_lasso <- sqrt(mean((y_pred-y_test)^2))
mae_lasso <- mean(abs(y_pred - y_test))
rsq_lasso <- cor(y_pred, y_test)^2


#code implementation for elastic net

#create the en regression with a warm start/ and find timing
#the alpha = 0.5 is the split between L1 and L2 regularization of the model
#remember that EN is the combination of the two!
 time_en <-system.time({
    en_model <- cv.glmnet(x,y, alpha=0.5, nfolds = 10,
                          lambda=seq(0.001, 0.1, length=100), standardize=TRUE,
                          keep=TRUE, type.measure="mse", warm.start=TRUE)


 })


 #now, use the model to make predictions
x_test <- model.matrix(medv~., data = test)[,-1]
y_test <- test$medv
y_pred <- predict(en_model, newx = x_test, s = en_model$lambda.min)


#calculate comparison criterion
```

```r
rmse_en <- sqrt(mean((y_pred-y_test)^2))
mae_en <- mean(abs(y_pred - y_test))
rsq_en <- cor(y_pred, y_test)^2
```

```python
# Note, I implemented this code in spyder but did not run it here


# -*- coding: utf-8 -*-
"""
Spyder Editor


This is a temporary script file.
"""


import pandas as pd
import numpy as np
from sklearn.linear_model import Ridge, Lasso, ElasticNet
from sklearn.datasets import fetch_openml
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from sklearn.model_selection import train_test_split, GridSearchCV
import time



boston = fetch_openml(name='boston', as_frame=True)


#Split the data
X, y = boston.data, boston.target


#need to fix the columns in x
X = X.apply(pd.to_numeric, errors='coerce')
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=4000769
```

```python
param_grid = {'alpha' : np.logspace(-4,4,9)}



#time the fitting for the warm start process/
start_time = time.time()


#fit the data initially for a starting point
#find optimal lambda


ridge = Ridge()


grid_search = GridSearchCV(estimator=ridge, param_grid = param_grid, cv = 10, scoring = "neg_me
grid_search.fit(X_train, y_train)


lambda_opt = grid_search.best_params_['alpha']


ridge_warm_started = Ridge(alpha = lambda_opt)



ridge_warm_started.fit(X_train, y_train)
end_time = time.time()



#Compute the metrics
y_pred = ridge_warm_started.predict(X_test)
rmse_ridge = mean_squared_error(y_test, y_pred, squared = False)
mae_ridge = mean_absolute_error(y_test, y_pred)
rsq_ridge = r2_score(y_test, y_pred)
time_ridge = end_time - start_time


#Repeat the process for Lasso
```

```python
lasso = Lasso(alpha = 1, warm_start=True)

start_time = time.time()

grid_search = GridSearchCV(estimator=lasso, param_grid = param_grid, cv = 10, scoring = "neg_me
grid_search.fit(X_train, y_train)

lambda_opt = grid_search.best_params_['alpha']




lasso = Lasso(alpha = lambda_opt, warm_start=True)
lasso.fit(X_train, y_train)

end_time = time.time()

#Compute the metrics
y_pred = lasso.predict(X_test)
rmse_lasso = mean_squared_error(y_test, y_pred, squared = False)
mae_lasso = mean_absolute_error(y_test, y_pred)
rsq_lasso = r2_score(y_test, y_pred)
time_lasso = end_time - start_time

#Repeat it again for Elastic Net

en = ElasticNet(alpha = 1)

start_time = time.time()

grid_search = GridSearchCV(estimator=en, param_grid = param_grid, cv = 10, scoring = "neg_mean_
grid_search.fit(X_train, y_train)
```

```python
lambda_opt = grid_search.best_params_['alpha']



en = ElasticNet(alpha = lambda_opt)
en.fit(X_train, y_train)


end_time = time.time()


#Compute the metrics
y_pred = en.predict(X_test)
rmse_en = mean_squared_error(y_test, y_pred, squared = False)
mae_en = mean_absolute_error(y_test, y_pred)
rsq_en = r2_score(y_test, y_pred)
time_en = end_time - start_time
```