

EE 239AS
Special Topics in Signals and Systems
Project 1
Regression Analysis
Winter 2016

Liqiang YU, Kaiming WANG and Jun FENG
904592975, 504592374, 304588434

01-31-2016

Contents

1	Introduction	3
2	Data Preprocessing	3
2.1	Feature Binarization	3
2.2	Pattern Detection	3
3	Linear and Random Forest Regression	4
3.1	Linear Regression	4
3.2	Random Forest Regression	6
3.3	The comparison between two models	6
4	Neural Network	8
5	Piece-wise Linear Regression and Polynomial Regression	9
5.1	Piece-wise Linear Regression	9
5.2	Polynomial Regression	9
6	Linear and Polynomial Regression of Boston Datasets	9
6.1	Linear Regression	9
6.2	Polynomial Regression	12
7	Lasso & Ridge Regularization	12
7.1	Lasso Regularization	13
7.2	Ridge Regularization	13

1 Introduction

In this report, we implemented some regression models to fit the data within two datasets, including the linear regression, random forest regression, neural network regression and polynomial regression. We made some comparisons between the prediction results of different models based on root mean squared error(RMSE) and the coefficient of determination(R^2). Based on the statistics provided by models, we analyzed the significance of different features. In order to avoid overfitting, which often happens in the regression fitting problem, we implemented both cross-validation and regularization techniques. Moreover, since some of the features have no numerical meaning, we formatted them into binary features, which increase the number of features.

The report is organized as follows : In section 2, we preprocessed the data and plotted the data to find some pattern. In section 3, The dataset was from a simulated traffic data on the backup system in a network system. We implemented both the linear regression model and random forest regression model to predict the copy size of the backup file and made some analysis and comparisons of the predictions. In section 4, we implemented the neural network model and polynomial regression model to fit the data and analyzed the impact of different parameters to the results. In section 5, we implement the piecewise linear regression and polynomial regressions. In section 6, we implement all the previous regression on the Boston Housing Dataset. In section 7, we introduce ridge and lasso regression techniques on this Dataset

2 Data Preprocessing

2.1 Feature Binarization

After checking the data, there are seven features in the dataset named network_backup_dataset, including week, day of the week, backup start time, work flow ID, file name, size of backup and backup time. The target value is the size of backup, others are features. Therefore, first of all we need to transform non-numerical values into numerical values. However, consider "the day of the week" feature, "Tuesday" is not larger than "Monday", but "2" is larger than "1". So we need to make such kind of features into binary format. For instance, for "day of the week", we need to expand it into seven features, where it is 1 on the specified day and all the others 0. Finally, we have 64 features in total.

2.2 Pattern Detection

In order to find the pattern, we plot the figure of the copy size over time. The independent variable is the "Start Time - Hour of Day", the dependent variable is the copy size. There are six sampling times in a day and 20 days as a period. So the x axis is from 1 to 120. The figure is shown in figure 1. From 1 we can see there is clearly a cycle between copy size and time. So we plot the copy size

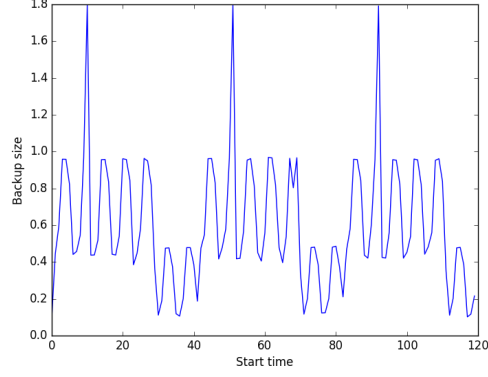


Figure 1: The copy size of backup over start time

over days to find the cycle period. The result is shown in figure 2. Hence we can predict that the cycle period is approximately a week.

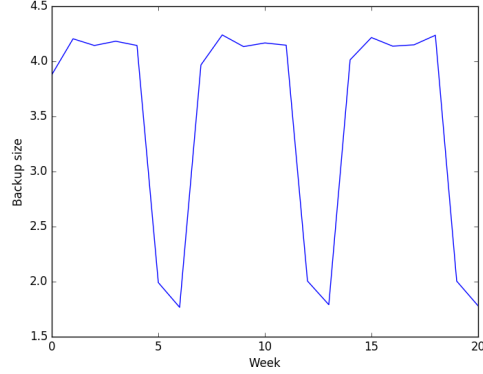


Figure 2: The copy size of backup over week

3 Linear and Random Forest Regression

3.1 Linear Regression

In the Linear Regression Model, we choose the copy size as the target variable and the other attributes as the features, the fit model is as follow

$$\bar{Y} = X\alpha$$

where \bar{Y} represents the target, X represents the features vector and α is the linear coefficient. We used least square as the penalty function, which is

$$\min ||Y - \bar{Y}||_2$$

In order to avoid the overfitting problem, we implemented 10-fold cross-validation technique. First do a random shuffling over the data, then split them into 10 folds. Each time choose 9 of 10 as the train set and the other 1 fold as the test set. The Root Square Mean Error (RMSE) is shown in figure 3. The average RMSE is 0.071. To evaluate the fitting accuracy of our model, Fitted values VS

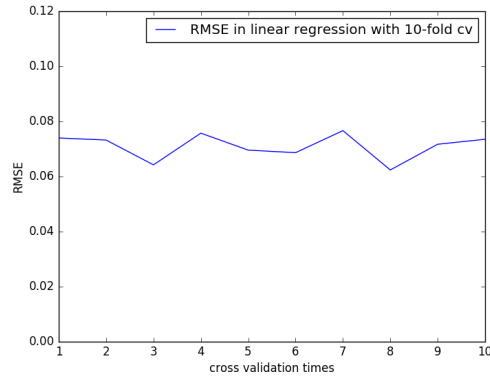


Figure 3: The copy size of backup over week

actual values and residuals vs fitted values are plotted. From figure 4, we can see that most predicted values are close to the actual values. From figure 3.1, we can see that the residuals are distributed randomly around zero axis, which means the model is proper.

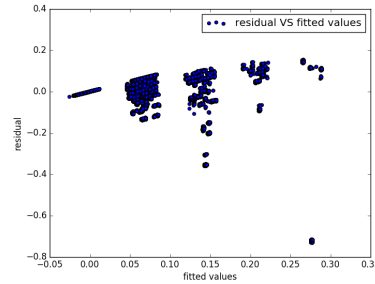
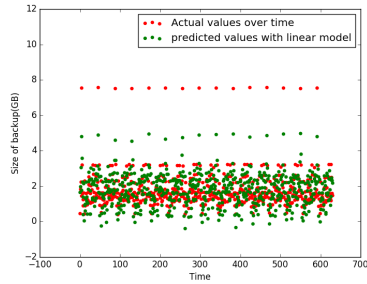


Figure 4: Actual VS predicted

Figure 5: Residual VS fitted values

As for the significance of different variables, we choose p-value as the evaluation

criterion. The p-value of the six features is shown in figure 6.

```
[ 7.78443589e-01  8.71062213e-04  2.46112227e-08  1.32748780e-01
 1.57542554e-01  0.00000000e+00]
```

Figure 6: p-value of 6 features

Therefore the most important two features are "Backup start time - Hour of the day" and "Backup Time" because their p-values are close to 0. For 64 features, the p-value matrix is shown below.

```
7.36199210e-001  9.39383338e-001  4.66439149e-001  9.08462632e-001
5.45241529e-001  4.35994229e-001  3.09252234e-001  6.38013695e-001
5.16578878e-001  2.80070834e-001  4.85768453e-001  8.55993609e-001
6.36535181e-001  7.63174401e-001  2.98120865e-001  1.18985257e-004
9.31487572e-004  1.66246279e-001  1.95819809e-001  4.44216789e-001
5.17062882e-001  1.99647310e-024  6.99156954e-009  1.54736893e-004
1.23481345e-002  3.53017379e-014  5.81362330e-002  1.24279127e-001
1.59934030e-033  8.85184242e-006  7.79480962e-034  2.34530310e-060
2.97238246e-100  2.68855252e-006  1.23242625e-005  3.15955319e-004
1.10012870e-005  5.78513377e-005  2.40988217e-007  2.36501171e-003
3.07991642e-001  4.75428019e-002  6.68776922e-002  5.09616461e-002
8.84143253e-001  6.74472700e-005  9.82830391e-006  6.18754039e-006
7.10107131e-003  9.41370915e-007  2.92999661e-010  5.36294108e-008
2.37168301e-008  1.50927326e-011  4.51686551e-009  2.17186539e-010
1.27307840e-007  6.34206490e-013  3.52571891e-010  4.98335413e-016
5.48046529e-013  1.95679737e-017  3.45472611e-014  0.00000000e+000]
```

Figure 7: p-value of 64 features

3.2 Random Forest Regression

For random forest regression, we initialize the model with parameters : number of trees 20, depth of trees 4, max features 64. The average RMSE is 0.0297. After tuning the parameters, we have the best RMSE 0.00943. The parameters are : number of trees 32, depth of trees 12, max features 64. The best RMSE is shown in figure 8.

3.3 The comparison between two models

The comparison between Linear Regression and Random Forest Regression is shown in figure 9. Random forest is a lot better than linear regression.

The prediction of the random forest model is shown in figure 10. We can see the cycle period is around 42 backup times, which is about a week, so the pattern is the same with the actual values.

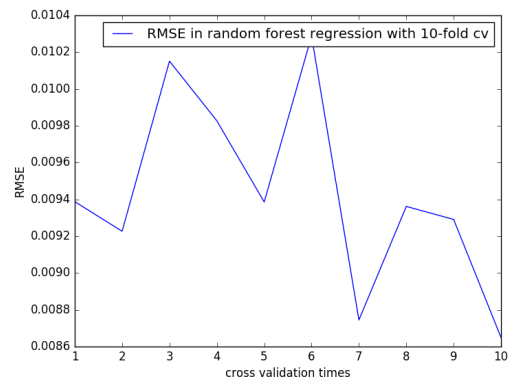


Figure 8: The best RMSE of random forest regression

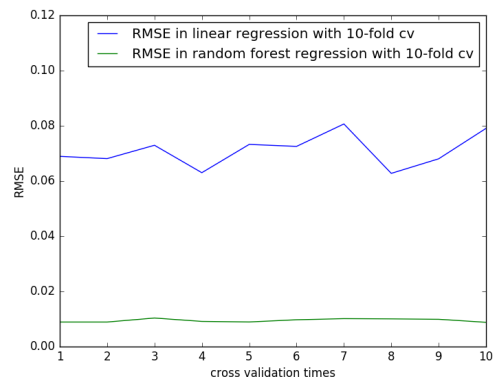


Figure 9: The comparison between Linear and Random Forest Regression

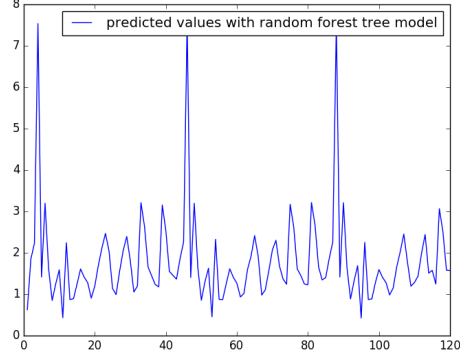


Figure 10: The prediction of the random forest

4 Neural Network

In the neural network model, we choose the feedforward network with the back-propagation trainer. The main parameters are the number of layers, the kind of hidden layers and the number of iteration before convergence. After lots of experience, we find that sigmoid layer is proper for hidden layers and linear layer is proper for the output layer. The optimal number of layers is 7. The RMSE reduces when the number of iteration is increasing, however it is too slow with large iteration number. So we set it to be 10. The best RMSE we can get is 0.0827. The RMSE is shown in figure 11.

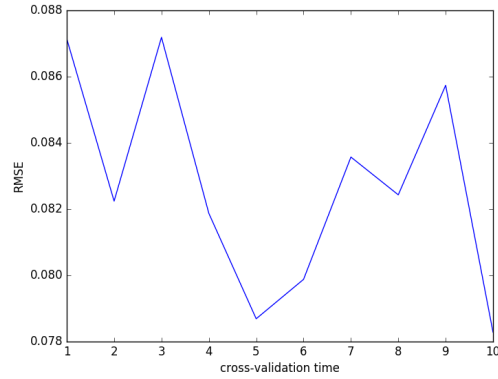


Figure 11: The RMSE of neural network model

5 Piece-wise Linear Regression and Polynomial Regression

5.1 Piece-wise Linear Regression

We just divide the all tuples according to their work-flow #, then implement the linear regression respectively which is illustrated by figure 12. However, we find that sometimes, piecewise regression does not necessarily benign the RMSE. It is consistent to the most cases: at most time, more data and more features can bring about a more reliable regression.

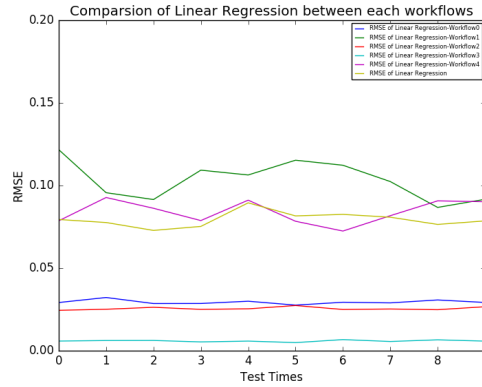


Figure 12: Piecewise Linear Regression

5.2 Polynomial Regression

By implying the PolynomialFeatures of Scikit-Learn, we can do polynomial regression. However, when the degree go higher and higher, overfitting may happened. Just as in figure 13. According to the figure, the threshold of the degree is 5. Cross validation can help us to decide the degree of polynomial function.

6 Linear and Polynomial Regression of Boston Datasets

6.1 Linear Regression

Similar as the linear part of Chapter 3, we also implemented 10-fold cross validation here. And the RMSEs are shown in Figure 14, and the mean of them is about 4.7403. To evaluate how well our model fit the data, we provide two graphs. The first one is shown as Figure 15, which is the scattered plot of actual

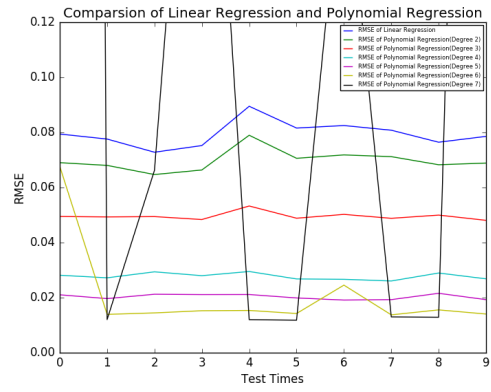


Figure 13: Polynomial Regression

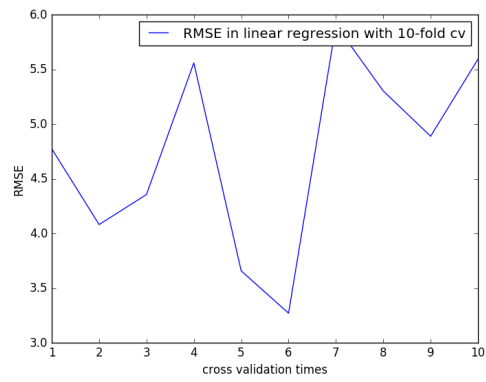


Figure 14: Linear Regression of Boston

and fitted MEDV values. And we can see that most of them overlaps. The other graph is the residuals versus fitted values, showed in Figure 16. Since most of them are clustered around 0 of y-axis, we are confirmed that our model fits the data quite well. For the analyzation of significance of different features, we

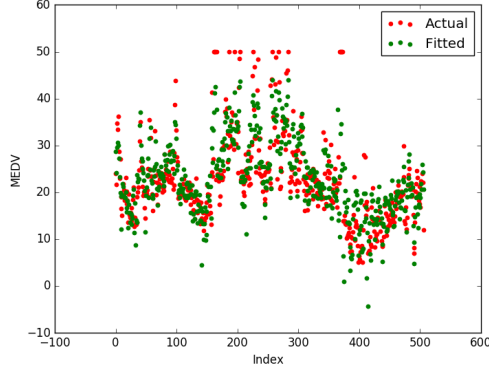


Figure 15: Scattered Plot of Actual and Fitted MEDV Values

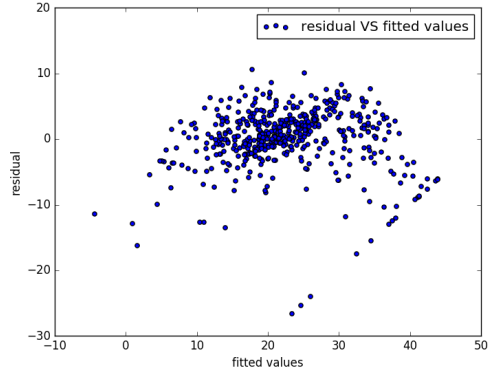


Figure 16: Residuals Versus Fitted Values

still use the p-value as before. P-value of the 13 features are shown in Figure 17, while those of 21 features are shown in Figure 18. From them, we know that the three smallest p-values are the sixth, the third last and the last, which means that the corresponding features are the most significant ones. Particularly, they are the RM, PTRATIO and LSTAT. Later, when we were using the Lasso regression model and the parameter is set as 0.1, we noticed that the only regression coefficients remain nonzero are exactly corresponding to those mentioned above, which means only the most important features are kept.

```
[ 5.50417920e-04  2.64479457e-01  6.12021303e-02  1.24477692e-02
 3.62538807e-02  1.47967137e-08  3.07657195e-02  7.27372323e-01
 3.31331513e-03  1.67962192e-03  1.18477865e-05  4.36407860e-03
 3.44786672e-15]
```

Figure 17: P-value of the 13 Features

```
[ 4.58294633e-02  5.81383430e-03  6.28136070e-07  3.32684686e-01
 1.51719694e-03  8.89434048e-12  4.39716459e-05  3.77815886e-02
 5.65818351e-01  3.36847735e-01  2.18034717e-01  4.14932326e-01
 2.55606463e-01  5.62375646e-01  9.57351116e-03  6.88691021e-02
 4.51512402e-04  1.86651257e-05  1.03757323e-03  3.01945047e-03
 1.08187699e-17]
```

Figure 18: P-value of the 21 Features

6.2 Polynomial Regression

We do the same polynomial regression on target of Boston dataset, just like in part 3, this time the optimal degree of fit is just 2, which is illustrated in Figure 19.

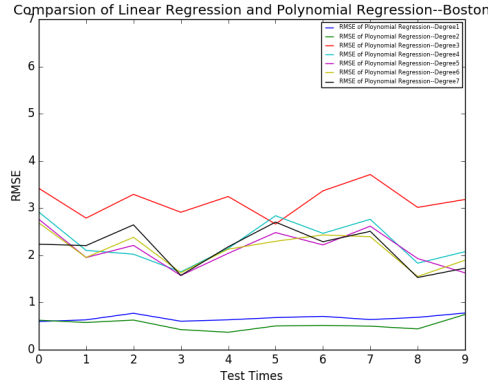


Figure 19: Polynomial Regression

7 Lasso & Ridge Regularization

In machine learning, over fitting is likely to happens when a model is excessively complicated. Under this situation, the machine only memorizes the training data, instead of learning to predict the data from given features. In our problem, we have only about 500 tuple of data with 20 different features. Thus, the chance of our model to memorize data from these features is relatively high. Before, we use cross validation to limit the affection of over fitting during the assessing of a model. In this way, we can only determine if a given model is good or not. Instead, if we want to generate better learning models with little over

fitting, we may use method as regularization. In regularization, we introduce a regularization term or so called loss function, which is usually a penalty on the complexity of our model. Therefore, it can improve the generalization ability of our model.

7.1 Lasso Regularization

In Lasso regularization, the regularization term is considering the l_1 -norm of the regression coefficient vector. And the penalty function has the following form:

$$\min \frac{1}{2n} \|Y - X\beta\|_2^2 + \alpha \|\beta\|_1$$

For different extent of penalty on the complexity, we can choose various value of α . And in our solution, we choose it in the range 0.1, 0.01, 0.001 as required. The corresponding performance of the regression model is shown in Figure 20, from which we can see that the RMSE is smaller with a little value of α . After further

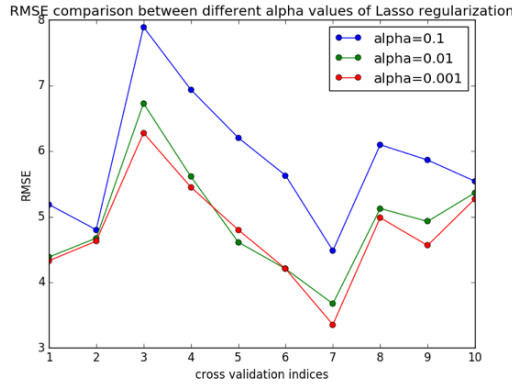


Figure 20: Lasso Regularization of Boston Dataset

research, we know that the result of linear regression almost overlaps the result of $\alpha = 0.001$, which means that the RMSE of pure linear regression is smaller than those of Lasso regression. Thus, in this problem, the control of over fitting does not bring a better model to us. However, this does not mean it is useless. After checking the result of the regression coefficients, we are shocked by its sparsity. When $\alpha = 0.1$, only 3 coefficients are nonzero and all the other 17 are 0. Thus, in this way, we have found a great sparse solution of our problem.

7.2 Ridge Regularization

In Ridge regularization, the regularization term is considering the l_2 -norm of regression coefficient vector. And the penalty function has the following form:

$$\min \frac{1}{2n} \|Y - X\beta\|_2^2 + \alpha \|\beta\|_2$$

Table 1: Compare between Lasso regression and Ridge regression

Value of α	Lasso regularization	Ridge regularization
0.1	5.8639	4.7981
0.01	4.9307	4.7686
0.001	4.7861	4.7725

Similar as Lasso regularization, we set the value of α in the range $\{0.1, 0.01, 0.001\}$ and the performance of our regression models are shown in Figure 21. This time, the curves of RMSE for different α are very close to each other. For a clearer

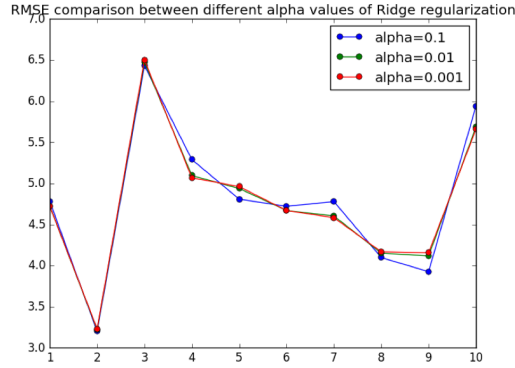


Figure 21: Ridge Regularization of Boston Dataset

comparison of the all RMSE values, we calculate the mean RMSE upon every 10-fold cross validation and we put all of them in Table 1. From the above table, we find out that for Lasso regularization, the RMSE is worsened but the sparsity of regression coefficients is obtained. And for Ridge regularization, we get the best RMSE result at $\alpha=0.01$, which prove that with the control of over fitting, we may get a better model.