

# Sparse Regression

Ankur Ankan, Kevin Jacobs, Zhuoran Liu

## 1 Introduction

### 1.1 Lasso Regression using Coordinate Descent

Lasso Regression is a regularized linear regression technique in which we use  $L_1$  regularization and get the following optimization problem:

$$\min_{\beta} \frac{1}{2} \sum_{\mu=1}^p \left( y^{\mu} - \sum_{i=1}^n \beta_i x_i^{\mu} \right)^2 \quad (1)$$

under the constraint that:  $\sum_{i=1}^n |\beta_i| \leq t$ . The iterative solution to this optimization problem using coordinate descent is given by:

$$\beta_j \leftarrow S \left( \frac{1}{p} \sum_{\mu} \tilde{y}_j^{\mu} x_j^{\mu}, \gamma \right) \quad (2)$$

where:

$$\begin{aligned} \tilde{y}_j^{\mu} &= y^{\mu} - \sum_{i \neq j} \beta_i x_i^{\mu} \\ S(\hat{\beta}, \gamma) &= \text{sign}(\hat{\beta})(|\hat{\beta}| - \gamma)_+ \end{aligned} \quad (3)$$

Lasso regression works particularly well in the cases when the number of features are huge. Unlike Ordinary Linear Regression, Lasso has the tendency to converge the weight values towards 0 and hence using only more important features and avoiding overfitting.

## 2 Sequential Gauss-Seidel rule for Lasso

Starting with the Langrangian form of Lasso optimization problem, we need to minimize  $f$  given by:

$$f = -\frac{1}{2p} \sum_{\mu=1}^p \left( y^{\mu} - \sum_{i=1}^n \beta_i x_i^{\mu} \right)^2 + \sum_{i=1}^n \gamma |\beta_i| \quad (4)$$

Taking the derivative with respect to  $\beta$  and equating it to 0 we get:

$$\frac{\partial f}{\partial \beta_j} = \frac{1}{p} \sum_{\mu=1}^p \left( y^\mu - \sum_{i=1}^n \beta_i x_i^\mu \right) x_j^\mu + \gamma \text{sign}(\beta_j) = 0 \quad (5)$$

Putting it in matrix form we have:

$$-b_j + \sum_i \chi_{ij} \beta_i + \gamma \text{sign}(\beta_j) \quad (6)$$

$$\chi \beta' = b - \gamma \text{sign}(\beta) \quad (7)$$

Now from Gauss-Seidel rule we have the iterative solution of the equation  $AX = b$  as:

$$x'_i = \frac{1}{A_{ii}} \left( b_i - \sum_{j>i} A_{ij} x_j - \sum_{j<i} A_{ij} x'_j \right) \quad (8)$$

Comparing the update rule with the equation we get:

$$\beta'_i = \frac{1}{\chi_{ii}} \left( (b - \gamma \text{sign}(\beta))_i - \sum_{j>i} \chi_{ij} \beta_j - \sum_{j<i} \chi_{ij} \beta_j \right) \quad (9)$$

### 3 Research Questions

1. How does the accuracy of the model change for different values of  $\gamma$  ?
2. How does  $\gamma$  affect the absolute value of  $\beta$  ?
3. How does Lasso Regression perform when the features are correlated ?

### 4 Results

We are given a dataset of 50 samples each having 100 features. The first step is to standardize our data set using:

$$X_j = (X_j - \text{mean}(X_j)) / \text{std}(X_j) \forall j = \{1, 2, \dots, n_{\text{features}}\} \quad (10)$$

We did a simple test for the fit of our learned parameter values shown in Fig. 1. We have plotted our predicted  $y$  against the most important feature i.e. the feature with the highest weight  $X[0]$ . We can see in the figure that even with a single feature we are able to predict the data quite well.

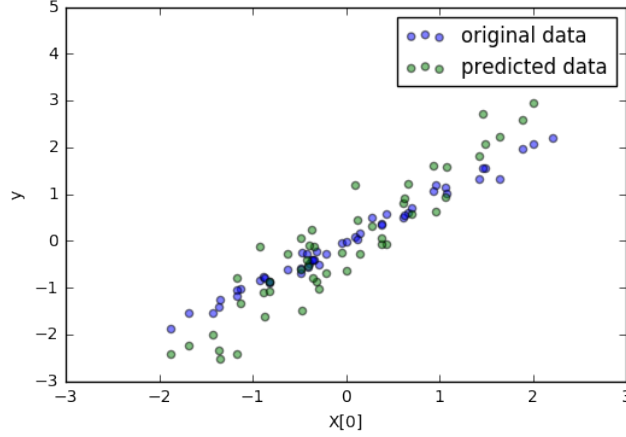


Figure 1: Original target variables with the predicted values against the 1st feature ( $X[0]$ ) of the data. Model learned with  $\gamma = 0.1$

#### 4.1 Variation of parameters with change in gamma

We start with random values for  $\beta$  sampled from a Normal Distribution and iteratively update each value of  $\beta$  one by one keeping the others fixed. In Fig. 2, we can see the variation of the weights with the iterations. We can see that the weights converge quite quickly in just 4 iterations and most of the weight values converge to 0 as we expect in the case of Lasso regression. Fig. 3 shows the number of parameters that converge to 0 for different values of  $\gamma$ . As we increase the value of  $\gamma$ , we penalize the weights more and hence more parameters start converging to 0 as we can see in the figure.

#### 4.2 Variation of accuracy with change in gamma

In Fig. 4 we have plotted the mean squared error on the validation set for different values of  $\gamma$ . We see that when  $\gamma$  is really small we have a huge value of  $t$  and hence most of the parameters have non-zero weights and hence it results in overfitting because of which we are getting high mean squared error. Also in the case when  $\gamma$  is high we have a too constrained model resulting in underfitting and hence the high mean squared error. We get the best accuracy on the validation set for  $\gamma = 0.22$ .

#### 4.3 Performance in the correlated case

For the correlated case we given a dataset having 1000 samples with 3 features each. From the correlation matrix in Fig. 6, we can see that the 2nd feature is correlated with both 0 and 1. It is more correlated with 1st feature than the 0th one. Therefore for good predictions regression should assign highest weight

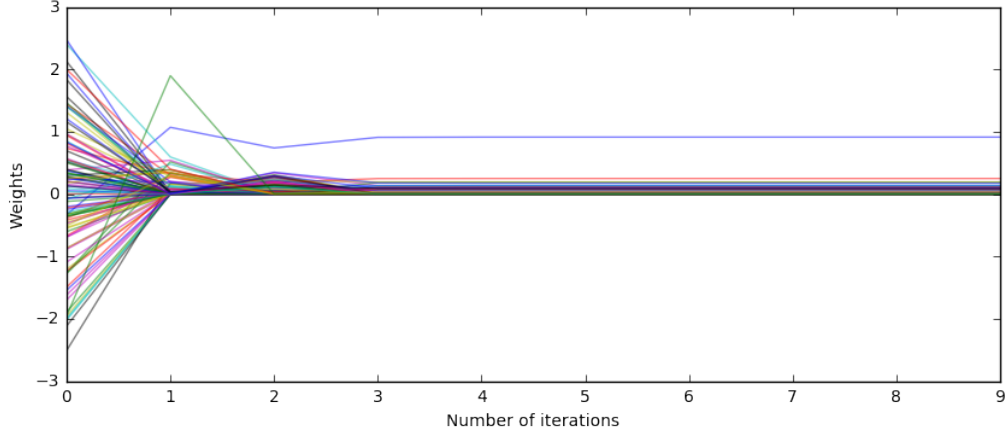


Figure 2: Change in weights  $\beta$  with iterations. In each iteration each  $\beta$  value is updated exactly once.

to the 1st feature and least weight to 2nd feature. But we can see in Fig. 7 that both 0th and 2nd features are turned off for lower values of  $\gamma$  and then the 2nd feature has the highest weight, resulting in the poor performance of the mode.

We also compared the performance of the Lasso Regression with Ridge Regression in Fig. 5. We can see that Ridge regression performs better than Lasso for all the values of gamma.

## 5 Conclusion and Discussion

The LASSO (Least Absolute Shrinkage and Selection Operator) is a regression method that involves penalizing the absolute size of the regression coefficients. By penalizing we end up in a situation where some of the parameter estimates may be exactly zero. The larger the penalty applied, the further estimates are shrunk towards zero. This is convenient when we want some automatic feature selection algorithm or when our feature size is too huge.

But as we saw one of the major drawbacks of this method is that it performs poorly in the case of correlated features and arbitrarily assigns 0 weight to some of the correlated case resulting in very poor performance. In such cases ridge regression is a good alternative as it is biased towards assigning a non zero weight to all the features.

## 6 Appendix

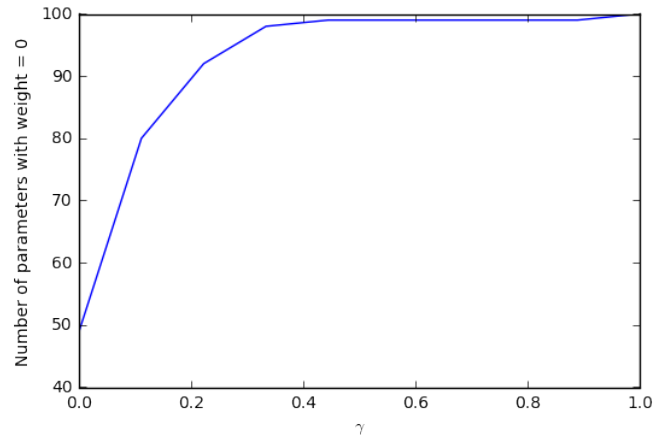


Figure 3: Number of parameters that have converged to 0 for various values of  $\gamma$

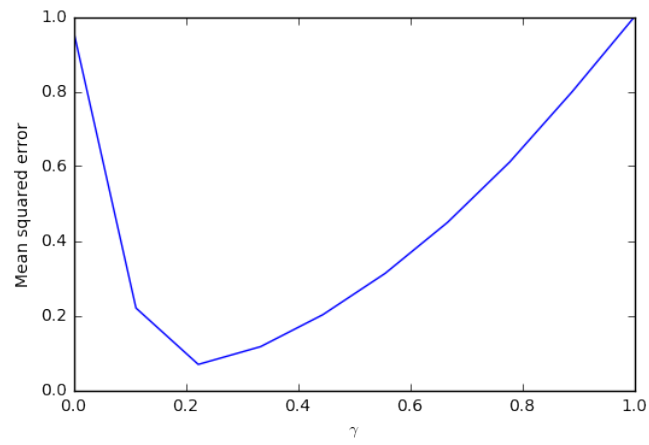


Figure 4: Mean Squared Error vs  $\gamma$  on validation dataset

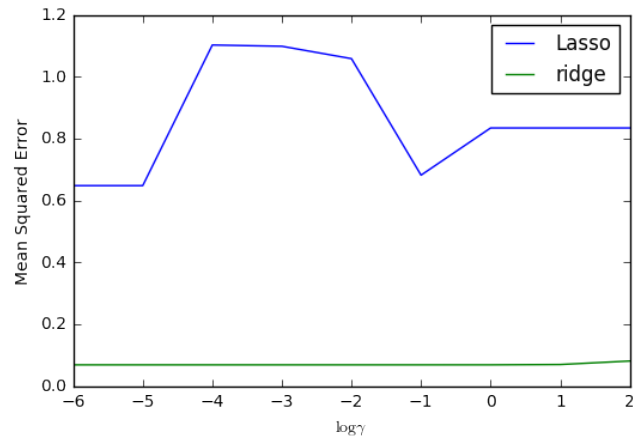


Figure 5: Comparison of Lasso and Ridge regression when the features are correlated

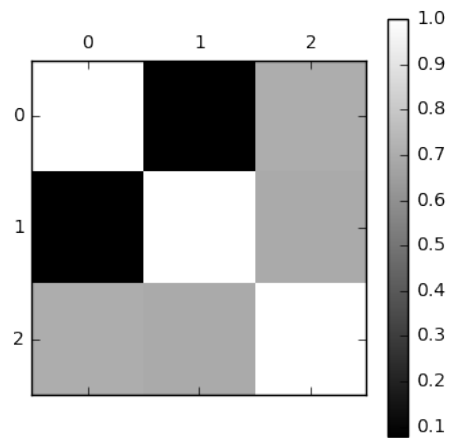


Figure 6: The correlation matrix between the 3 features of the given dataset

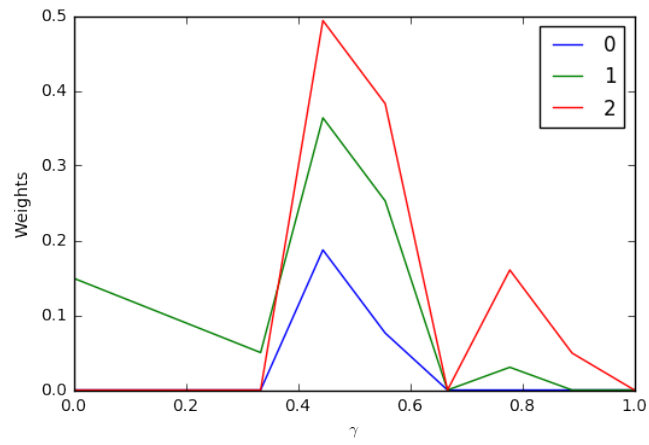


Figure 7: Converged values of  $\beta$  for different values of  $\gamma$  in the correlated dataset