

# Machine Learning 6.867 - Pset 1

October 1, 2015

## 1 Implement Gradient Descent

### 1.1 Implementation

We implemented a basic gradient descent function in MATLAB. In this function, the user can specify the objective function and the corresponding gradient function. The user also inputs an initial guess, the constant step size, and the convergence threshold. The algorithm works as follows: starting at the initial guess, we move to a new point by taking a step (of the specified step size) in the direction of the gradient. We repeat the process until the difference in objective value between two subsequent iterations is less than the specified threshold.

### 1.2 Testing Gradient Descent

We test the gradient descent method on the following functions:

1.  $f(x_1, x_2) = 3x_1^2 + 2x_1x_2 + x_2^2 - 4x_1 + 5x_2$

By completing the square we find that the optimal solution is:  $x_1 = 2.25, x_2 = -4.75$ , and the optimal value is:  $f(x_1, x_2) = -16.375$ .

2.  $f(x) = -\frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$

This is the negative of a Gaussian pdf, with  $\mu = 0$  and  $\sigma = 1$ . The optimal solution is  $x = 0$  and the optimal value is:  $f(x) = -\frac{1}{\sqrt{2\pi}}$ .

3.  $f(x) = x^3 - 10x$ .

This is not a convex function. There is one local minimum:  $x = 1.8257, f(x) = -12.17$ , but the global minimum is negative infinity.

The basic gradient descent method works well on the first two functions. For most specifications of the initial guess and step size, the procedure converges to the right answer. Increasing the convergence threshold to very large values leads to imprecise results, and decreasing it to very small values results in more iterations until convergence. In some cases, very large or very small step sizes may result in non-convergence. For the third function which is non-convex, selecting an initial guess  $x_0 \geq -1.8257$  leads us to the incorrect, local minimum.

### 1.3 Numerical Gradient

We can also numerically evaluate the gradient at a given point. The partial derivative at  $x_i$  is approximated by the formula below, where  $\mathbf{h}$  is the vector of all zeros except for the  $i$ th component that is  $\epsilon$ . We apply this numerical gradient in place of the analytic one in the original three functions. For sufficiently small  $\epsilon$ , the results are almost identical.

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = \frac{f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x} - \mathbf{h})}{2\epsilon} \quad (1)$$

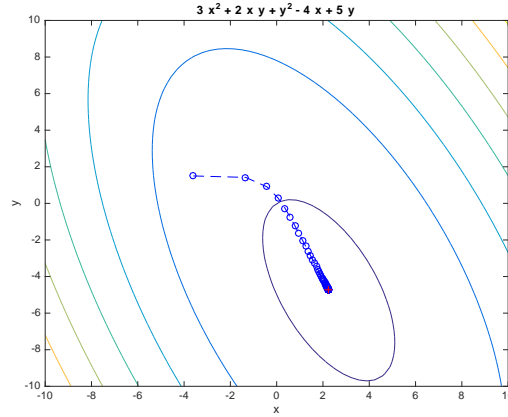


Figure 1: Contour plot of the first function. The gradient descent method moves in steps (blue dots) towards the local/global optimum (red cross).

## 1.4 Comparing with Sophisticated Methods

As the final step to evaluate our implementation of gradient descent, we compare our method against state-of-the-art gradient descent packages in MATLAB. Compared against `fminunc`, our implementation is dramatically slower. For example, given the same convergence criterion, it takes 35 iterations for our method to converge, while `fminunc` converges after 8 iterations. We think the disadvantage is due to the constant step size in our method. Since the step size does not adapt while the algorithm is running, our method is likely to undershoot or overshoot the optimal solution. An improvement could be made using backtracking, exact line-search, or other techniques to adaptively update the step size.

# 2 Linear Basis Function Regression

## 2.1 MLE for Polynomial Basis Function Regression

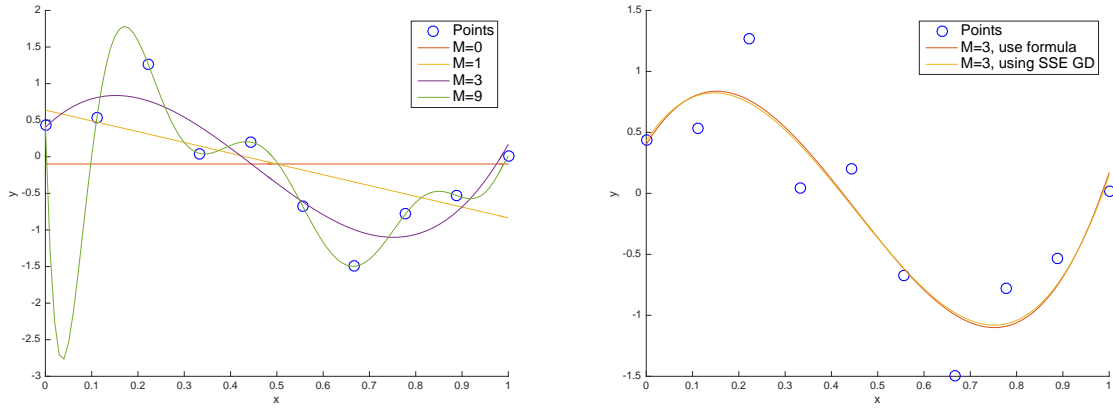
We implemented a maximum likelihood estimator for a linear regression with polynomial basis functions, using the formula  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ . The replicated weight vectors are close to what is reported in Table 1.1, and the plots we generated in Figure 2a match the ones in Bishop well.

## 2.2 Using SSE

Alternatively, the maximum likelihood weight vector can be found using numerical gradient descent on the SSE. We first implement code to calculate the SSE:  $(\mathbf{Y} - \Phi \mathbf{w})^T (\mathbf{Y} - \Phi \mathbf{w})$ . The derivative is  $-2\Phi(\mathbf{Y} - \Phi \mathbf{w})$ , which can also be computed numerically using the finite differencing method implemented in Problem 1. We found that the two methods yield identical gradient values.

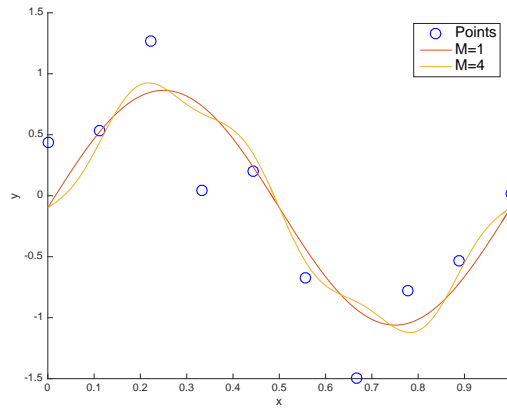
## 2.3 Use Gradient Descent on SSE to find MLE

To compute the MLE for  $\mathbf{w}$ , we input the SSE as the objective function in the gradient descent method, with  $\mathbf{w}$  as the variable to be optimized over. With a smart initial guess, step size, and reasonably small convergence threshold, our  $\mathbf{w}$  from gradient descent matches the analytic OLS solution closely (within a small margin of error). For example, if the initial guess is within 5% of the optimum, step size: 0.05, and convergence threshold:  $10^{-6}$ , then we attain solutions within 1% of the true weight vector  $\mathbf{w}$ . The comparison can be found in Figure 2b. Some bad choice of step size results in non-convergence or extremely large estimates. We also compare the speed of the method against the native MATLAB function `fminunc`. Given the same convergence criterion, it takes 42 iterations for our method to converge, while `fminunc` converges after 29 iterations.



(a) Replication of Bishop 1.4, using polynomial basis function up to degree  $M$ .

(b) Comparing the estimates using MLE equation or numeric gradient descent on SSE.



(c) Using sine basis function.

Figure 2

## 2.4 Sine Basis Function

We can also use sine basis functions, with  $\phi_1(x) = \sin(2\pi x), \dots, \phi_M(x) = \sin(2\pi Mx)$ . The estimated plots are in Figure 2c. We observe that the predicted functions are periodic as expected. One potential disadvantage is that when the underlying data is aperiodic, this estimation cannot capture the trend of non-periodic behavior. This drawback may be worse for regression problems with out-of-sample data that lies outside the range of observed features, because periodic functions have finite range.

## 3 Ridge Regression

### 3.1 Implementation

Ridge regression is the particular case of regularized least squares with a quadratic regularizer term. The error function that we aim to minimize over is given by:

$$\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad (2)$$

The closed-form solution of this problem is well-known, and can be derived by setting the gradient of (2) equal to zero. The optimal solution for  $\mathbf{w}$  is provided by Bishop (2006), page 145:

$$\mathbf{w}_{ridge} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (3)$$

We coded this method in MATLAB and tested our program using data from Bishop Figure 1.4, varying the parameters of  $\lambda$  and  $M$ . For the extreme cases, we observed that if  $\lambda \leq 0.0001$ , then  $\mathbf{w}_{ridge} \approx \mathbf{w}_{OLS}$ , and if  $\lambda \geq 100$ , then  $\mathbf{w}_{ridge} \approx \mathbf{0}$ . If  $M = 1$ , then the ridge regression prediction is linear, and as  $M$  increases, our prediction curves to fit the training data more closely. Predictions resulting from different choices of the parameters  $\lambda$  and  $M$  are shown in Figure 3.

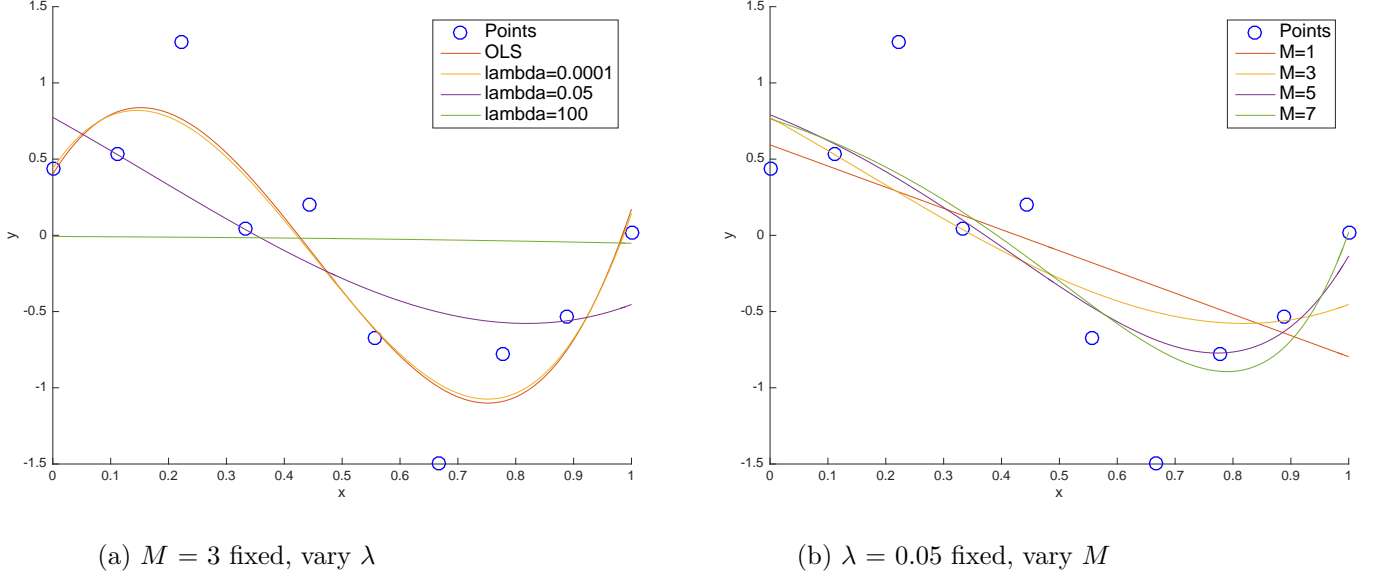


Figure 3: Prediction curves from ridge regression using data from Bishop 1.4.  $M$  = order of polynomial basis function,  $\lambda$  = regularization parameter.

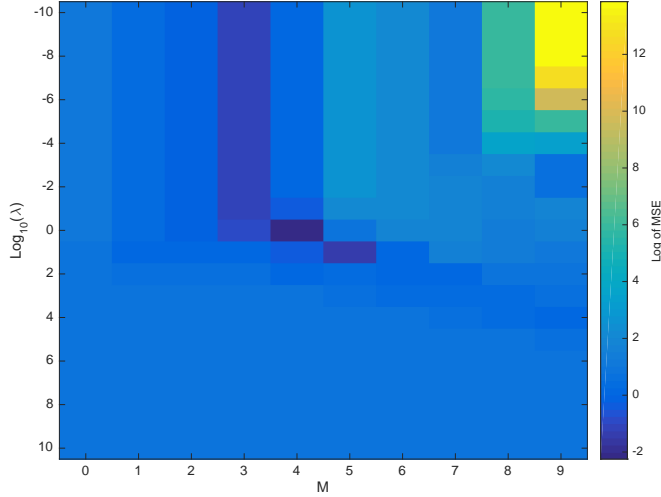
### 3.2 Model Selection

To optimize parameter values for  $\lambda$  and  $M$ , we build our models using training data and then compare out-of-sample performance on validation data. For this example, we performed a grid search over the ranges:  $\log_{10}(\lambda) = \{-10, \dots, 10\}$ ,  $M = \{0, \dots, 9\}$ . We found that the model with  $\lambda = 1$  and  $M = 4$  yields the lowest MSE on validation data, so we select these to be the final parameter values. This model leads to  $\text{MSE} = 0.1056$  on the validation data and  $\text{MSE} = 2.6988$  on the test data.

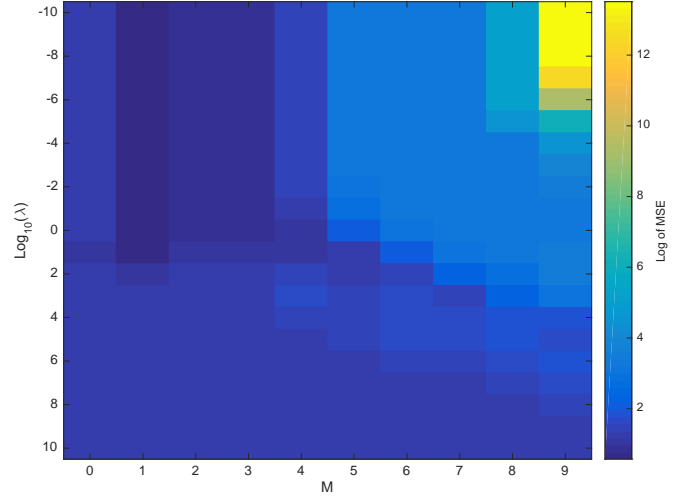
In general, we observe that models with relatively low MSE on the validation set tend to have relatively low MSE on the test set. Heat maps of the MSE values for different parameter values are shown in Figure 4. As expected, test MSE is higher on average, and the plots are similar for most parameter values. Discrepancies between the heat maps may be attributed to the presence of an outlier in the training set. In both heat maps, we see that the top right corner (high  $M$ , low  $\lambda$ ), which corresponds to complex models with little regularization, lead to the worst performance on out-of-sample data.

### 3.3 BlogFeedback Dataset

To test our ridge regression algorithm in a real-world application, we apply our method to the BlogFeedback regression problem from Krisztian Buza at the Budapest University. In the dataset, a blog post is described by 280 features, and our goal is to predict the number of comments it will receive in the next 24 hours. We are given a training, validation, and test set of blog post observations, and assume a simple linear basis function for this problem.



(a) Validation MSE



(b) Test MSE

Figure 4: Heat map of the MSE for ridge regression on the validation and test sets, varying the parameters  $M$  and  $\lambda$ . Color gradient from low MSE (navy blue) to high MSE (yellow).

Doing a line search over the range  $\log_{10}(\lambda) = \{-4, \dots, 10\}$ , we find that  $\lambda = 10,000$  yields the lowest MSE on validation data. This model leads to  $\text{MSE} = 1014.4$  on the validation data and  $\text{MSE} = 894.36$  on the test data.

## 4 Sparsity and LASSO

An alternative approach is to use the  $L_1$  norm regularizer in the objective for regularized least squares. The error function that we aim to minimize over is given by:

$$\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j| \quad (4)$$

The LASSO is widely used in practice because it leads to solutions which are relatively sparse, i.e. the resulting prediction vector  $\mathbf{w}$  tends to have few non-zero components.

To test this method and its sparsity properties, we consider a regression problem with 5 training points and 500 testing points in  $\mathbb{R}^{12}$ . The data is generated according to  $f(x) = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon$ , where  $\epsilon \sim N(0, \sigma^2)$  and  $\mathbf{w} \in \mathbb{R}^{12}$  is a sparse vector with only two non-zero components. For this problem, we assume that  $\mathbf{w}$  is unknown, and we estimate this vector with ridge regression in Section 4.1 and with LASSO in Section 4.2. Let  $\lambda_2$  denote the regularizer term for the  $L_2$  norm and let  $\lambda_1$  denote the regularizer term for the  $L_1$  norm. Figure 5 shows plots of the estimated function with different regularizers.

### 4.1 Ridge Regression Gradient Descent

Applying our method for gradient descent to the ridge regression problem on the training dataset of 12 points, we compute the vector of weights  $\mathbf{w}$  and MSE on the testing dataset for  $\lambda_2 = 0.1$ . We obtain:

$$\begin{aligned} \mathbf{w} &= (0.12, 0.14, 0.22, 0.20, 0.09, -0.05, -0.20, -0.30, -0.32, -0.25, -0.08, 0.12) \\ \text{MSE} &= 0.2606. \end{aligned}$$

If we disable the  $L_2$  regularizer (set  $\lambda_2 = 0$ ), i.e., the OLS estimator, we obtain:

$$\mathbf{w} = (1.15, 1.30, 1.72, 0.67, -1.73, -4.28, -5.16, -3.27, 0.48, 3.01, 0.61, -8.83),$$

$$\text{MSE} = 109.2910.$$

## 4.2 LASSO Gradient Descent

Similarly, we apply our method for gradient descent to the LASSO problem on the on the training dataset of 12 points. For  $\lambda_1 = 0.1$ , we find:

$$\mathbf{w} = (0.00, 0.00, 0.43, 0.00, 0.00, -0.00, -0.00, -0.00, -1.02, -0.00, -0.00, 0.00),$$

$$\text{MSE} = 0.0487.$$

## 4.3 Sparsity Properties of LASSO

We observe from Figure 5 and MSE values that the LASSO estimation is closest to the true generating function, compared to ridge regression or non-regularized OLS. In terms of sparsity, the LASSO estimator has only two variables that are non-zero (the rest are within numerically negligible deviation from 0), whereas the other estimators lead to non-sparse results. The  $L_0$  norm regularizer, which penalizes based on the number of non-zero elements in  $\mathbf{w}$ , is the true sparsity regularizer; however, this leads to a problem which is non-convex and difficult to solve directly. Because the  $L_1$  norm is the convex hull of the  $L_0$  norm, LASSO can be used as an approximation to often give sparse results. Even though the  $L_2$  norm regularization penalizes the size of  $\mathbf{w}$ , due to the geometry of the norm it does not provide sparse results.

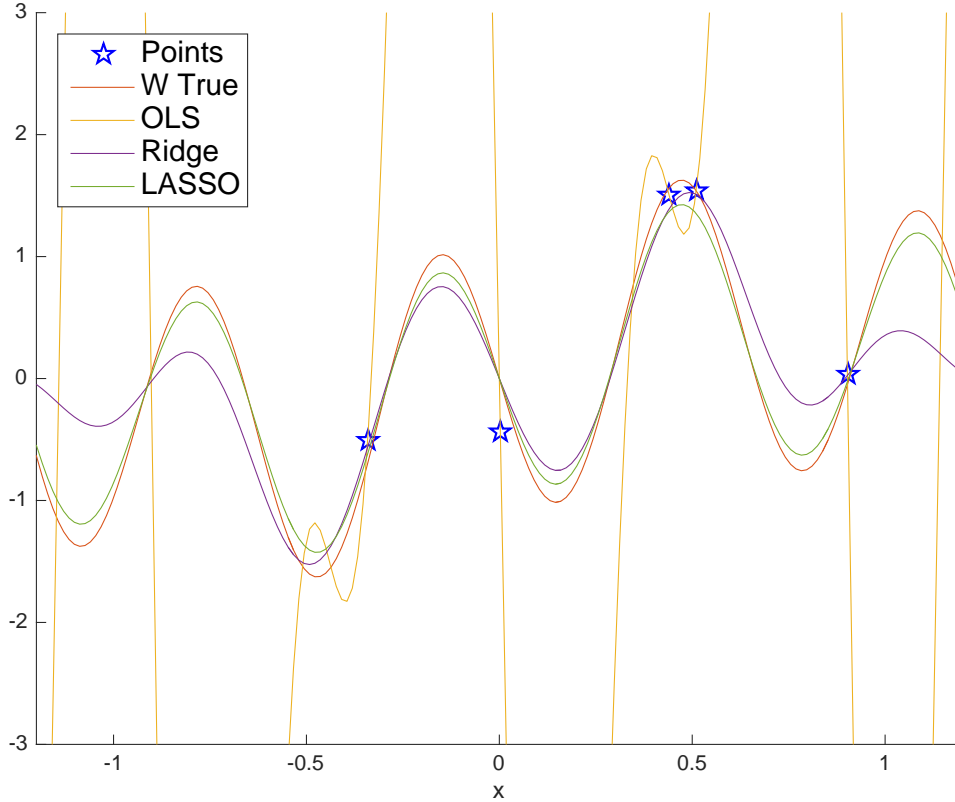


Figure 5: Plots of estimated function with different regularizers.  
 OLS:  $\lambda_1 = \lambda_2 = 0$ , Ridge:  $\lambda_2 = 0.1$ , LASSO:  $\lambda_1 = 0.1$ .