

# Machine Learning 6.867 - Pset 1

September 27, 2015

# 1 Implement Gradient Descent

## 1.1 Gradient Descent 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. User also specify initial guess, the constant step size, and the convergence threshold.

The procedure works as the following: 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 objective values between two iterations are less than the specified threshold.

## 1.2 Testing Gradient Descent

We test the gradient descent procedure in 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 get the optimal value is at  $x_1 = 2.25, x_2 = -4.75$ .)
2.  $f(x) = -\frac{1}{\sqrt{2\pi}\sigma}e^{-(x-\mu)^2/2\sigma^2}$  (negative of a Gaussian pdf, with  $\mu = 0$  and  $\sigma = 1$ .)
3.  $f(x) = x^3 - 10x$ . (This is not a convex function; local minimum is at 1.8257 and global minimum is negative infinity).

The basic gradient descent method works decently well in the first two functions. In most specifications of the initial guess and step size the procedure converges to the right answer. Changing the convergence threshold to very large gives imprecise results, and changing it to very small results in more iterations until convergence. In some cases, very large or very small step sizes may result in non-convergence. In the third function which is non-convex, the initial guess can lead us to the incorrect, local minimum.

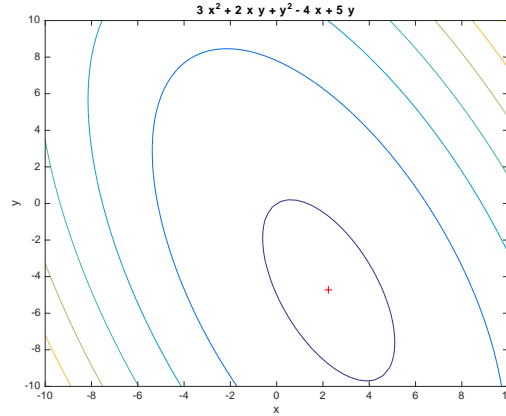


Figure 1: Contour plot of the first function. The gradient descent method moves to the red cross where the local/global optimal is.

## 1.3 Numerical Gradient

We can also numerically evaluate the gradient at a given point. The partial derivative at  $x_i$  is evaluated by taking the difference between  $f(\mathbf{x} + \mathbf{h})$  and  $f(\mathbf{x} - \mathbf{h})$ , where  $\mathbf{h}$  is the vector of all zeros except for the  $i$ th component being  $\epsilon$ , and divide the difference by  $2\epsilon$ . We apply this numerical gradient in place of the analytic one in the original three functions, and found that for sufficiently small  $\epsilon$ , the results are almost identical!

## 1.4 Comparing with Sophisticated Methods

As the final step to evaluate our implementation of gradient descent, we would like to compare ours with the more developed methods. Comparing with `fminunc`, our implementation is dramatically slower, unfortunately. For example, for the same convergence criterion, it takes 8 iterations for `fminunc` to converge, whereas it takes 35 iterations for our method. We think the disadvantage is from our constant step size. Since this size does not adapt to where we are, our method is likely to undershoot or overshoot. An improvement could be using backtracking, exact line-search, etc. methods to adaptively update the step size.

## 2 Ridge Regression

### 2.1 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:

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

The closed-form solution of this problem is well-known, and can be derived by setting the gradient of (1) equal to zero. The optimal solution for  $\mathbf{w}$  is provided by Bishop (2006), page 145:

$$\mathbf{w}_{ridge} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (2)$$

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  $\mathbf{w}_{ridge} \approx \mathbf{0}$ .

### 2.2 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:  $\lambda = \{0.0001, 0.001, 0.01, 0.1, 1, 10, 100\}$ ,  $M = \{0, 1, 2, 3, 4, 5\}$ . We found that the model with  $\lambda = \text{TODO}$  and  $M = \text{TODO}$  yields the lowest MSE on validation data, so we select these to be the final parameter values. This model leads to MSE = TODO for the test data, which is the TODO lowest overall. In general, we observe that models with MSE for the validation set tend to yield low MSE on the test set.

### 2.3 Blog Feedback