

Machine Learning Homework #1

Joe Ellis

September 26, 2013

1 PROBLEM 1

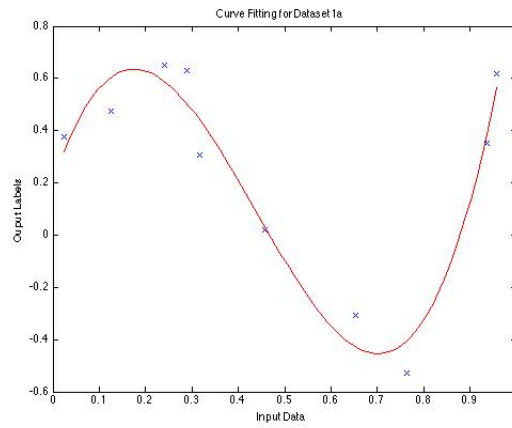
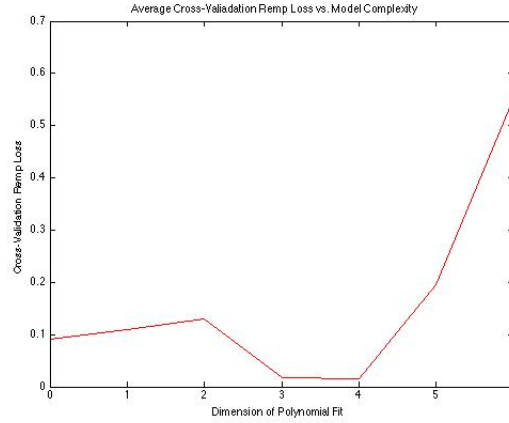
In this problem we investigate fitting a polynomial model and regression to a dataset consisting of 10 data points. Our d -dimensional polynomial function, f , is described below,

$$f(x; \theta) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d \quad (1.1)$$

where θ is the learned parameters of our model, and x is a given data point. We choose to use squared error for our empirical loss function to gauge how well our model fits the data present. The equation for empirical loss can be seen below, where x represents the given data points, y represents the output, which our function f is attempting to solve for.

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^n (y_i - f(x_i; \theta))^2 \quad (1.2)$$

We perform cross-validation to find the adequate dimensionality of an appropriate model. The cross-validation performed is completed 1000 times for each model complexity with train/test split of 9 points for training and 1 point for testing. The results of the cross-validation can be seen in Figure 1., and we can see that a model of 4 or 5 dimensions is suitable for this dataset. Models with lower than 3-dimensions have a hard time describing the variation in the data, whereas models with many more than 3-dimensions tend to overfit.



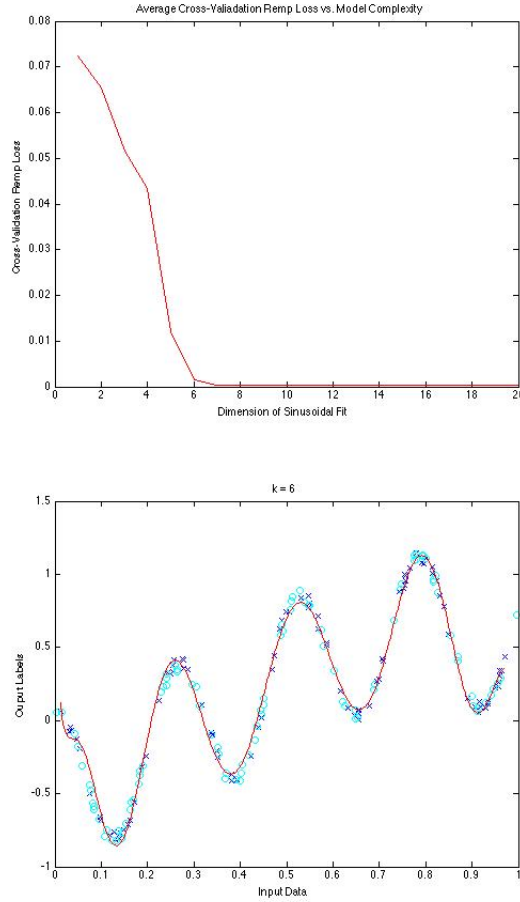
Using the 3-dimensional model for f we achieve the best polynomial regression function with with our $R_{emp}(\theta) = 0.0046$, and $\theta = (0.1951, 5.4884, -19.6041, 14.9130)$. The plot of the polynomial that was found to fit the data set can be seen in Figure 2.

2 PROBLEM 2

In Problem 2 we also address the issue of finding some function to fit to a given set of data points and their real valued labels. However, instead of using a standard polynomial basis function as described in Problem 1 we instead use a set of basis functions that consist of harmonic sets of *cos* and *sin* functions,

$$f(x; \theta) = \theta_0 + \sum_{i=1}^k \theta_i \sin(i * x) + \sum_{i=1}^k \theta_{i+k} \cos(i * x), \quad (2.1)$$

and these are known as sinusoidal basis functions. The calculation of the sinusoidal regression function shown above was completed in the “sinusoidalreg.” program, which



was handed in with this assignment. The $R_{emp}(\theta)$ was calculated in the same way as above. For this problem we once again used 1-dimensional input data with real valued output variables from “dataset1b.txt” that was provided to us. For this dataset the first 100 elements from dataset1b were used for training and the second 100 were used for testing. We completed 100-fold cross validation across the training set using a training/test split of 80 points for training and 20 for test to determine what was a good complexity for to model these data points. A plot of the cross-validation accuracies can be seen below in Figure 3.

We can see that a good dimensionality for k that allows for a simple model but good performance, is $k = 6$. Using $k = 6$ as our model parameter we were able able to achieve a training error of $R_{emptrain} = 0.0012$ and testing error of $R_{emptest} = 0.0042$. The found model has 13 parameters, and these parameters will be outputted to the screen if “runProblem2.m” is executed in the matlab command line environment, they were ommitted here for brevity. Finally, the result of our sinusoidal regression function can be seen below in Figure 4, where points marked by an “x” are training points, and those marked by an “o” are test points.

3 PROBLEM 3

In Problem 3 we will use a linear logistic regression function to perform a binary classification task on 2-D input data. In this problem we find a suitable θ for our logistic classification function by using batch gradient descent. The logistic classification function is defined as,

$$f(x; \theta) = \frac{1}{1 + \exp(-\theta^t x)}. \quad (3.1)$$

For this problem instead of using the average squared error for our loss function instead we will use, the logistic loss cost function, which is defined as,

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(x; \theta)) - y_i \log(f(x; \theta)). \quad (3.2)$$

To perform gradient descent to find the global minimum of $R_{emp}(\theta)$ we need to find the partial derivatives of $R_{emp}(\theta)$ with respect to each different parameter within the θ vector. The partial with respect to each element of θ is,

$$\frac{\partial R_{emp}}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x; \theta)) x_{ij}. \quad (3.3)$$

Once we have the partial derivative for each element we can find the $\nabla R_{emp}(\theta)$, and then update our theta using this found gradient. We define a stepsize as $\eta = 100$, and then update our θ after every iteration of gradient descent according to the equation,

$$\theta^{t+1} = \theta^t + \eta * \nabla R_{emp}(\theta) \quad (3.4)$$

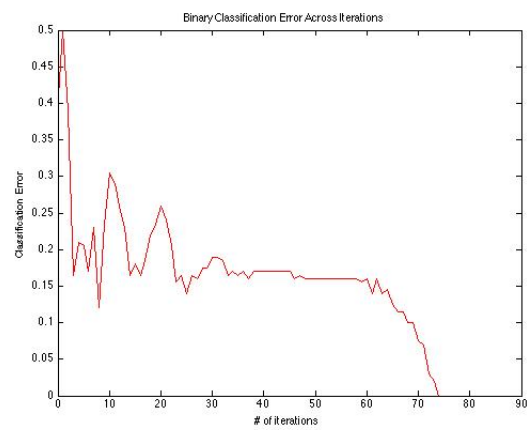
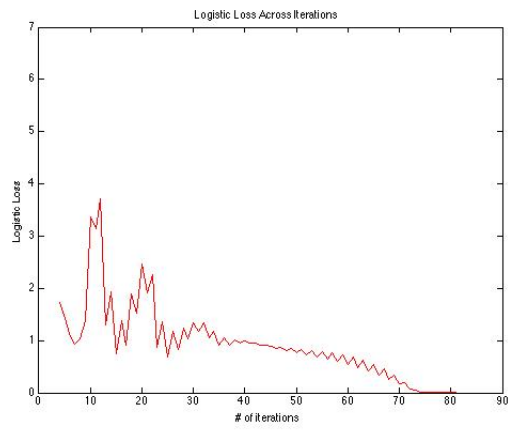
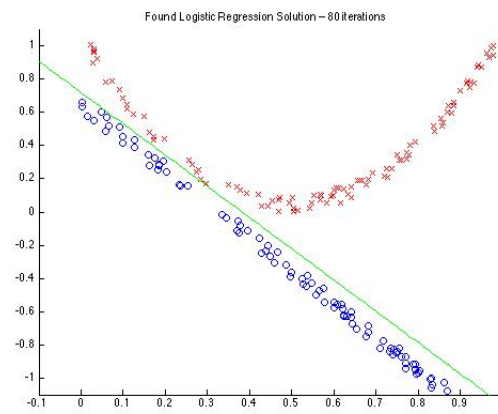
This algorithm completes once $|\theta^{t+1} - \theta^t| < \epsilon$, where we have chosen $\epsilon = 0.1$. Using the gradient descent algorithm we were able to achieve a perfect binary classification error on the given dataset, and our algorithm converges to a solution after only 80 iterations. The η chosen allowed for very aggressive steps, and with this our Loss and binary error had some level of oscillation throughout the process. Our final found normalized parameters give $\theta = (-0.319, 0.836, 0.444)$, with a $R_{emp}(x; \theta) = 0.0146$. In Figure 6, we can see the final linear logistic classifier imposed on top of our dataset, which shows the classification scheme.

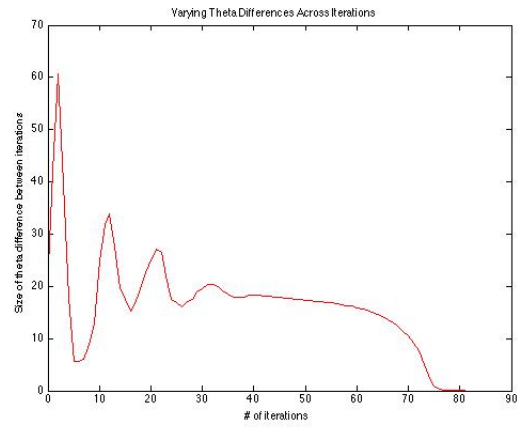
Figure 7 demonstrates how the $R_{emp}(\theta)$ varies over the iteration process, notice the oscillations due to the large step size.

Figure 8 shows how the Binary Classification error changes throughout the process.

The $|\theta^{t+1} - \theta^t|$ measure can be seen for every step in the iteration in Figure 9.

It should be noted that if a smaller η was chosen the plots 7-9 would have less oscillation within them. It is also true that we could possibly find a “better” solution if we allowed for smaller ϵ , but it would take longer to converge. A trade-off between speed and finding the best possible solution does exist with gradient descent.





4 PROBLEM 4