

# Comp 6321 - Machine Learning - Assignment 2

Federico O'Reilly Regueiro

October 18th, 2016

## Question 1:

### 1.a Partition the data into training / testing, 90 to 10 and perform and plot $L2$ regularization

The data was partitioned pseudo-randomly<sup>1</sup>.

```
phi = load('hw2x.dat');
phi = [phi, ones(size(phi,1),1)];
y = load('hw2y.dat');
```

```
% Partition the data randomly.
idxs = randperm(size(phi, 1));
idx_train = idxs(1:89);
idx_test = idxs(90:99);
```

```
phi_train = phi(idx_train, :);
y_train = y(idx_train);
phi_test = phi(idx_test, :);
y_test = y(idx_test);
```

Next, a range of lambdas was chosen, going from 0 to almost 14000 in order to get a good sense of the trend of both the error and the coefficients. The former was plotted in both a range of small values as well as along the whole set of lambdas for which the  $RMS$  error was calculated. This in order to allow us to see the behaviour of the test and training errors for lower values of  $\lambda$  as well as the overall trend of the  $RMS$  the resulting plot can be seen in Figure 1.

```
lambdas = 0:0.1:24;
lambdas = lambdas .^ 3;
```

```
w = zeros(length(lambdas), size(phi,2));
```

---

<sup>1</sup>I have included the permutation indexes yielded by matlab for the instance of the 90 / 10 partition from which the plots and values were drawn. Suffice it to uncomment two lines and comment-out two others in order to partition the data randomly, yet similar results (at different scales of  $\lambda$ ) can be observed.

```

for lambda = lambdas
    idx = lambda == lambdas;

    % Train the model
    w(idx, :) = pinv(phi_train' * phi_train ...
        + (lambda * eye(size(phi_train, 2)))) ...
        * (phi_train' * y_train);

    h_phi_train(:, idx) = phi_train * w(idx, :)';
    j_h_train(idx) = rms(h_phi_train(:, idx) - y_train);

    % Now compare to the test
    h_phi_test(:, idx) = phi_test * w(idx, :)';
    j_h_test(idx) = rms(h_phi_test(:, idx) - y_test);
end

```

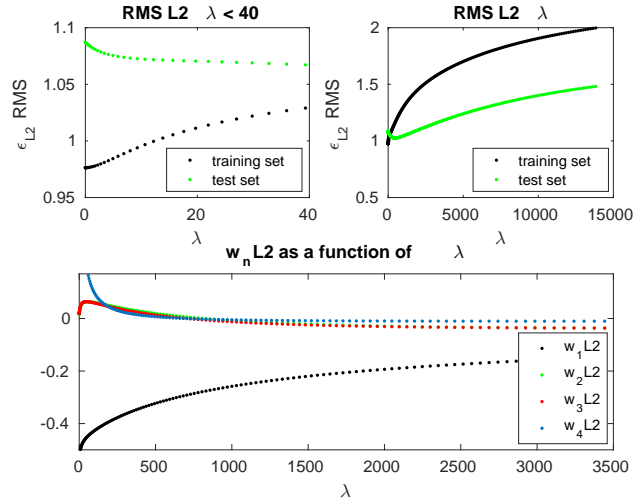


Figure 1: Plot of the RMS training and testing error as well as the coefficients over a wide range of  $\lambda$ .

On the top-left plot of figure 1 can observe that as expected, the test error goes down as  $\lambda$  grows. However, as  $\lambda$  continues to increase, the test error closely follows the training error as they both grow since the restrictions on the coefficients make for a worse fit at a certain point.

## 1.b Use the quadprog function in Matlab for $L1$ regularization

In order to use `quadprog` for regularization, we must first find the Hessian matrix  $\mathbf{H}$  as well as other parameters  $\mathbf{f}$ ,  $\mathbf{A}$ ,  $\mathbf{b}$  in the specific format that Matlab requires.

We recall that for  $L1$  regularization the expression we must minimize is:

$$\arg \min_w \frac{1}{2}(\Phi \mathbf{w} - \mathbf{y})^T(\Phi \mathbf{w} - \mathbf{y}) + \frac{\lambda}{2} \sum_{k=0}^{K-1} |w_k| \quad (1)$$

Which is equivalent to finding:

$$\arg \min_w (\Phi \mathbf{w} - \mathbf{y})^T(\Phi \mathbf{w} - \mathbf{y}) + \lambda \sum_{k=0}^{K-1} |w_k| \quad (2)$$

And expands to:

$$\arg \min_w \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{y}^T \Phi \mathbf{w} + \mathbf{y}^T \mathbf{y} + \lambda \sum_{k=0}^{K-1} |w_k| \quad (3)$$

And for which we can remove the constant term  $\mathbf{y}^T \mathbf{y}$ , yielding:

$$\arg \min_w \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{y}^T \Phi \mathbf{w} + \lambda \sum_{k=0}^{K-1} |w_k| \quad (4)$$

Matlab's `quadprog`( $\mathbf{H}$ ,  $\mathbf{f}$ ,  $\mathbf{A}$ ,  $\mathbf{b}$ ) function, gives the optimal  $\mathbf{x}$  corresponding to the expression  $\arg \min_x \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{f}^T \mathbf{x}$ , subject to constraints  $\mathbf{A} \mathbf{x} \leq \mathbf{b}$ . We can thus take  $\mathbf{H} := 2\Phi^T \Phi$ , then  $\mathbf{f} := -2\mathbf{y}^T \Phi$ ,  $\mathbf{A} := \lambda \mathbf{P}$ , where for a system with  $n$  variables,  $\mathbf{P}$  is the matrix with  $2^n$  permutations of  $[b_1, b_2, \dots, b_n, 0]^T$ ,  $b \in \{-1, 1\}$  and lastly  $\mathbf{b} := c \vec{\mathbf{1}}$ , where  $\vec{\mathbf{1}}$  is an all-one vector of length  $2^n$  that places an upper bound to the expression  $\lambda \sum_{k=0}^{K-1} |w_k|$ , such that  $\sum_{k=0}^{K-1} |w_k| \leq \frac{c}{\lambda}$ . If we set  $c = 1$ , then we parametrize  $L1$  regularization simply by adjusting  $\lambda$  accordingly.

Thus we end up with the following code:

```
for lambda = lambdas
    idx = lambda == lambdas;
    w_quad(idx,:) = quadprog(2*(phi_train' * phi_train), ...
        -2*(phi_train' * y_train), ...
        lambda*[ 1, 1, 1, 0; 1, 1,-1, 0; ...
            1,-1, 1, 0; 1,-1,-1, 0; ...
            -1, 1, 1, 0; -1, 1,-1, 0; ...
            1,-1, 1, 0; -1,-1,-1, 0], ...
        [1; 1; 1; 1; 1; 1; 1; 1]);
```

---

<sup>2</sup>the last zero entry is to avoid the regularization of the intercept term

```

h_phi_train_quad(:, idx) = phi_train * w_quad(idx, :)';
j_h_train_quad(idx) = rms(h_phi_train_quad(:, idx) - y_train);

% Now compare to the test
h_phi_test_quad(:, idx) = phi_test * w_quad(idx, :)';
j_h_test_quad(idx) = rms(h_phi_test_quad(:, idx) - y_test);
end

```

### 1.c Plot $L1$ RMS and coefficients, $w$ against $\lambda$ and comment

In Figure 2, we can again notice how the test error slightly decreases for the lowest values of  $\lambda$  and then monotonically increases as the coefficients are all forced towards 0. We also note that  $L1$  yields a lower minimum error (1.0098) than  $L2$  regularization (1.0257).

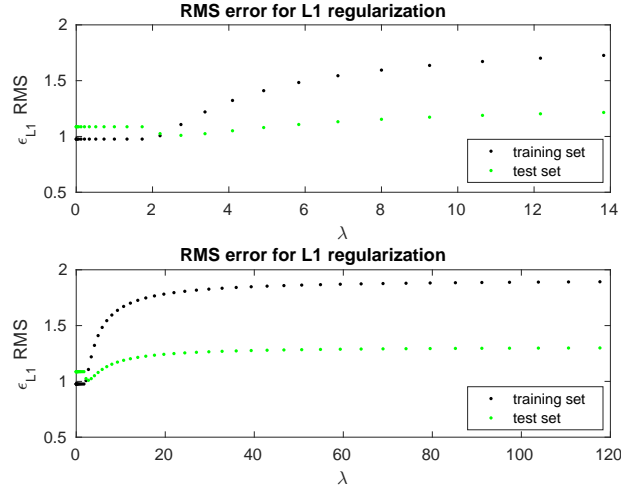


Figure 2: The RMS training and testing error for a wide range of  $\lambda$ .

By the same token, in figure 3 we can see how both  $w_2$  and  $w_3$  sharply decrease to 0 when  $\lambda = 2$  and the model relies solely on  $w_1$ <sup>3</sup> which then decreases gradually, as opposed to figure 1, where we can observe how all coefficients approach 0 at a similar rate during  $L2$  regularization. Conversely, as is expected we can see that both errors and coefficients are equal between  $L1$  and  $L2$  regularization when  $\lambda = 0$ .

<sup>3</sup> $w_4$  is the bias term, so we can't really say the model relies on it as it is not an input.

This particular data-set would lead to believe that the data was generated mainly by some function  $f(w_1)+\epsilon$ . This hypothesis is supported by the following output:

```
corr(y, phi(:,1:3))
ans =
    -0.848189    -0.017339    -0.024943
```

which reveals that the correlation between  $\Phi_{:,1}$  and  $y$  is much larger than between other columns of  $\Phi$  and  $y$ .

TODO change indexing from 1-based to 0-based

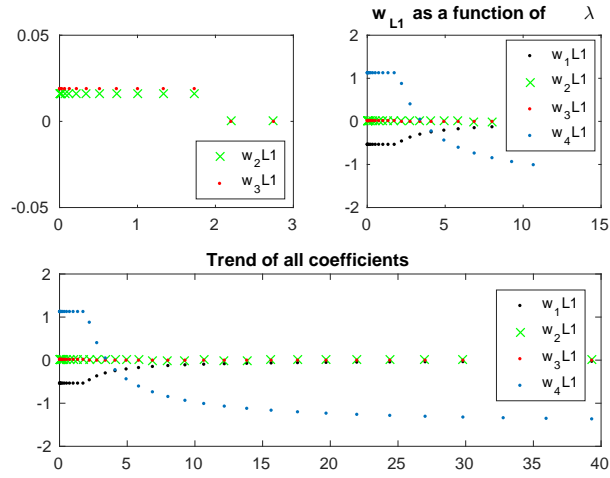


Figure 3: The RMS training and testing error for a wide range of  $\lambda$ .