

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-randomly¹.

```
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 λ 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));
```

¹I 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 λ) 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

```

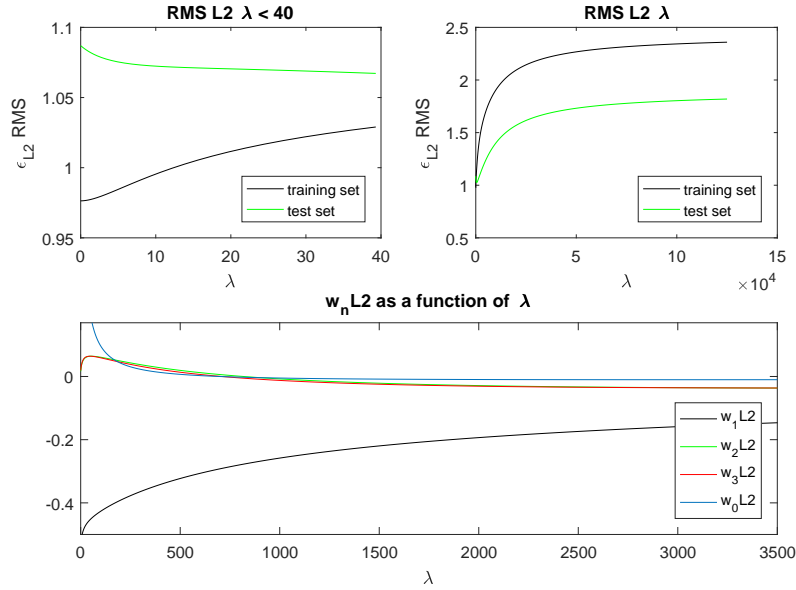


Figure 1: Plot of the RMS training and testing error as well as the coefficients over a wide range of λ .

On the top-left plot of figure 1 can observe that as expected, the test error goes down as λ grows. However, as λ 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 λ 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]);
```

²the 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 λ and comment

In Figure 2, we can again notice how the test error slightly decreases for the lowest values of λ 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).

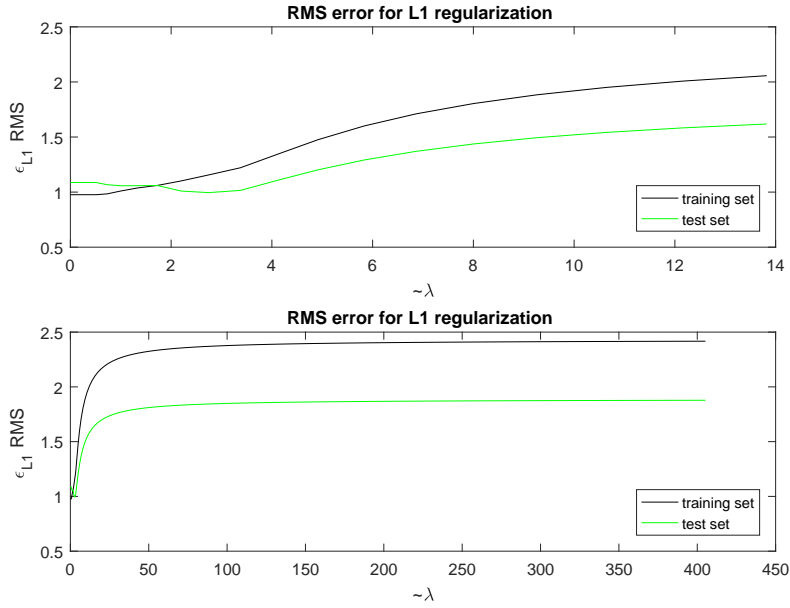


Figure 2: The RMS training and testing error for a wide range of λ .

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 ³ 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 $\lambda = 0$.

³ 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 Φ and y .

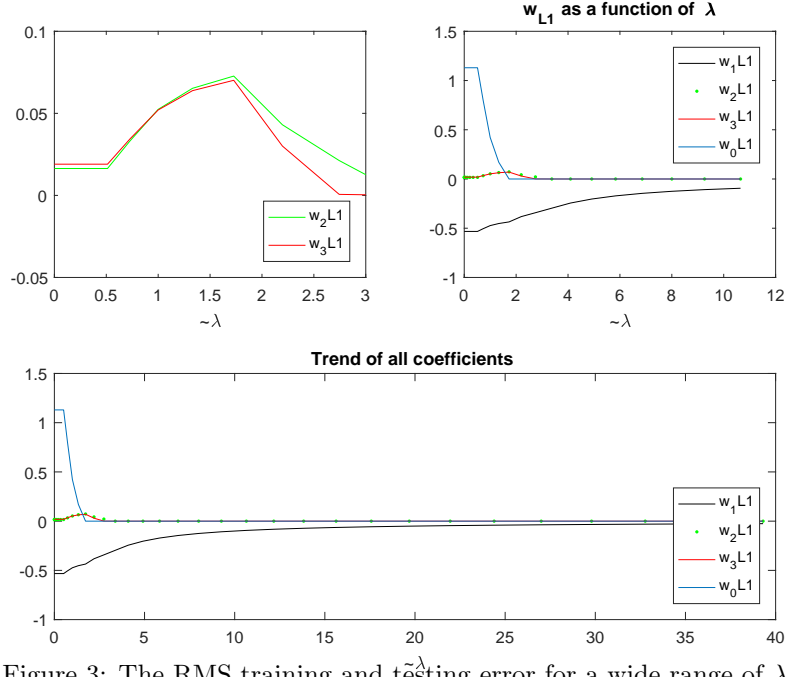


Figure 3: The RMS training and testing error for a wide range of λ .

iiiiiii HEAD

Question 2: Dealing with missing data, fill in $x_{i,n}$ with class-conditional means?

First, we write $\mu_{c,i}$ to represent $E(x_i|y = c)$ and we assume independence between features. Then, since our classifier is Gaussian, we know that $P(x|y = 1)$ and $P(x|y = 0)$ are modelled as follows:

$$P(x|y = c) = \frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(x - \mu_c)^T \Sigma^{-1}(x - \mu_c)}, \quad c \in \{0, 1\} \quad (5)$$

We turn our attention to the numerator of the exponent, which can also be written in the following manner:

$$\sum_{i=1}^n [(x_i - \mu_{c,i} \sum_{j=1}^n (\Sigma^{-1}_{i,j} (x_j - \mu_{c,j}))] \quad (6)$$

This expression becomes 0 . Thus we can write:

$$\log \frac{P(y=1|x)}{P(y=0|x)} = \log \frac{P(y=1)}{P(y=0)} + \log \frac{\frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_1)}}{\frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_0)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_0)}} \quad (7)$$

But since the matrix sigma is shared, we can then write:

$$\log \frac{P(y=1|x)}{P(y=0|x)} = \log \frac{P(y=1)}{P(y=0)} + \log \frac{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_1)}{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_0)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_0)} \quad (8)$$

Then we expand the exponent:

$$\log \frac{P(y=1|x)}{P(y=0|x)} = \log \frac{P(y=1)}{P(y=0)} + \log \frac{-\frac{1}{2} \sum_{i=1}^n [(x_i - \mu_{c,i} \sum_{j=1}^n (\Sigma^{-1}_{i,j} (x_j - \mu_{c,j}))]}{-\frac{1}{2} \sum_{i=1}^n [(x_i - \mu_{c,i} \sum_{j=1}^n (\Sigma^{-1}_{i,j} (x_j - \mu_{c,j}))]} \quad (9)$$

Where we can clearly see that the contribution of x_n does not change the ratio of the log-odds given by all other features, since we have choosen $x_n = \mu_{c,n}$ for $c \in 0, 1$ where $\mu_{c,n}$ is the class-conditional means for class c

Question 3: Naive Bayes assumption, suppose a feature gets repeated in the model