

Comp 6321 - Machine Learning - Assignment 2

Federico O'Reilly Regueiro

October 26th, 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 125000 in order to get a good sense of the trend of both the error and the coefficients. The former was plotted for a range of small values as well as along the whole set of lambdas for which the RMS error was calculated. Two plots were done 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.

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 pseudo-randomly, yet similar results (at different scales of λ) can be observed.

```

lambdas = 0:0.1:50;
lambdas = lambdas .^ 3;

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

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

```

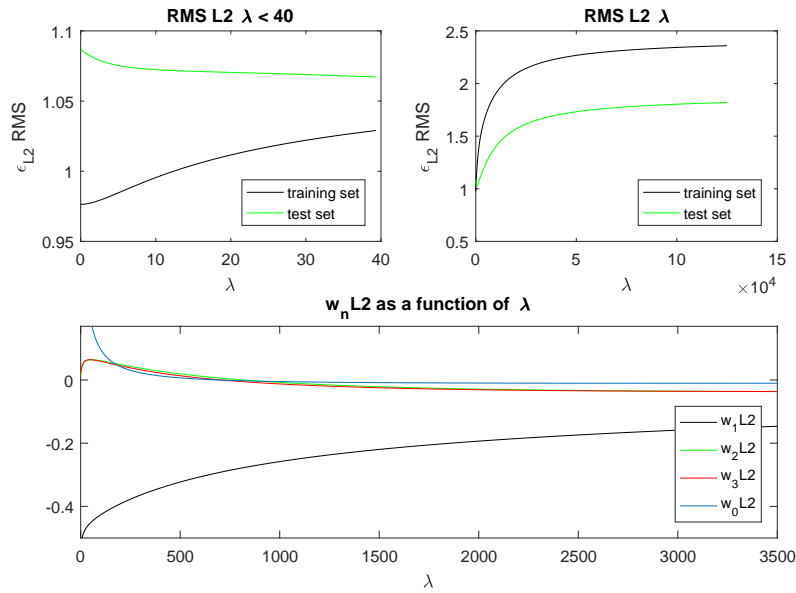


Figure 1: Plot of the RMS training and testing error as well as the coefficients over a wide range of λ .

On the top-left plot of figure 1, as expected, it can be noted that 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 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:

$$\begin{aligned} \arg \min_w (\Phi \mathbf{w} - \mathbf{y})^T(\Phi \mathbf{w} - \mathbf{y}) \\ \sum_{k=0}^{K-1} |w_k| \leq \eta \end{aligned} \quad (2)$$

And expands to:

$$\begin{aligned} \arg \min_w \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{y}^T \Phi \mathbf{w} \\ \sum_{k=0}^{K-1} |w_k| \leq \eta \end{aligned} \quad (3)$$

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

$$\begin{aligned} \arg \min_w \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{y}^T \Phi \mathbf{w} \\ \sum_{k=0}^{K-1} |w_k| \leq \eta \end{aligned} \quad (4)$$

Matlab's `quadprog(H, f, A, 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 $\mathbf{f} := -2\mathbf{y}^T \Phi$, $\mathbf{A} := \mathbf{P}$, where for a system with n variables, \mathbf{P} is the matrix with 2^n permutations of $[b_1, b_2, \dots, b_n]$, $b \in \{-1, 1\}$ and lastly $\mathbf{b} := \eta \vec{\mathbf{1}}$, where $\vec{\mathbf{1}}$ is an all-one vector of length 2^n that places an upper bound, η , to the expression $\sum_{k=0}^{K-1} |w_k|$, such that $\sum_{k=0}^{K-1} |w_k| \leq \eta$. We note that η is roughly equivalent to $\frac{1}{\lambda}$ which we use so that we may compare the effects of $L1$ and $L2$ regularizations on a similar range of values.

Thus we end up with the following code:

```

etas = 1./lambdas;
for eta = etas
    idx = eta == etas;
    w_quad(idx,:) =
        quadprog(
            2*(phi_train'*phi_train), ...
            -2*(phi_train'*y_train), ...
            [ 1, 1, 1, 1; 1, 1,-1, 1; ...
              1,-1, 1, 1; 1,-1,-1, 1; ...
              -1, 1, 1, 1; -1, 1,-1, 1; ...
              1,-1, 1, 1; -1,-1,-1, 1; ...
              1, 1, 1,-1; 1, 1,-1,-1; ...
              1,-1, 1,-1; 1,-1,-1,-1; ...
              -1, 1, 1,-1; -1, 1,-1,-1; ...
              1,-1, 1,-1; -1,-1,-1,-1], ...
            eta * [1; 1; 1;-1; 1; 1; 1; 1]);

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

    % Now validate
    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

Although it is not exactly equivalent, we shall simplify the comparison of $L1$ and $L2$ throughout this section by using $\lambda \approx \frac{1}{\eta}$ as it gives a clearer idea. In Figure 2, we can again notice how the test error slightly decreases for the lowest values of $\approx \lambda$ 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).

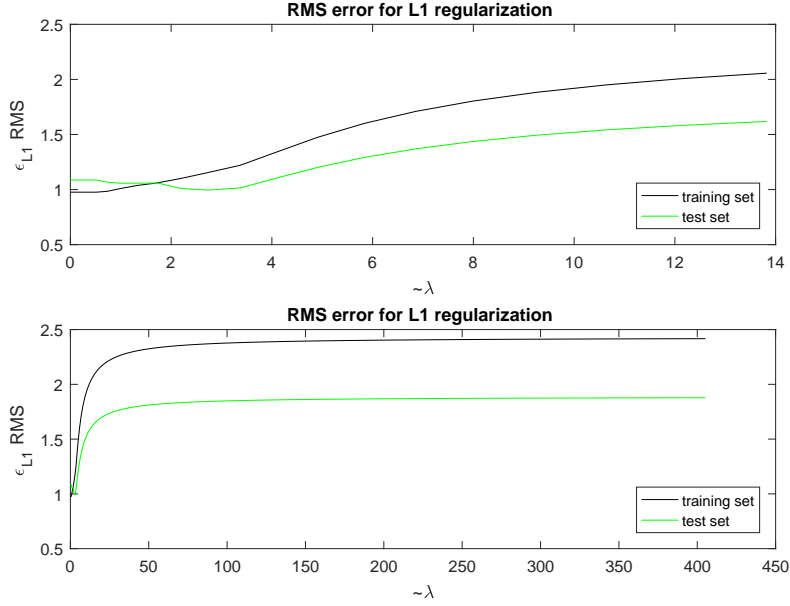


Figure 2: The RMS training and testing error for a wide range of $\frac{1}{\eta} \approx \lambda$.

By the same token, in figure 3 we can see how both w_2 and w_3 sharply decrease to 0 when $\approx \lambda = 2$ and the model relies solely on w_1 ² 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 and as expected, we can see that both errors and coefficients are equal between $L1$ and $L2$ regularization when $\lambda = 0$.

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 cross-correlation between $\Phi_{:,1}$ and y is much larger than between other columns of Φ and y .

² w_0 is the bias term, so we can't really say the model relies on it as it is not an input.

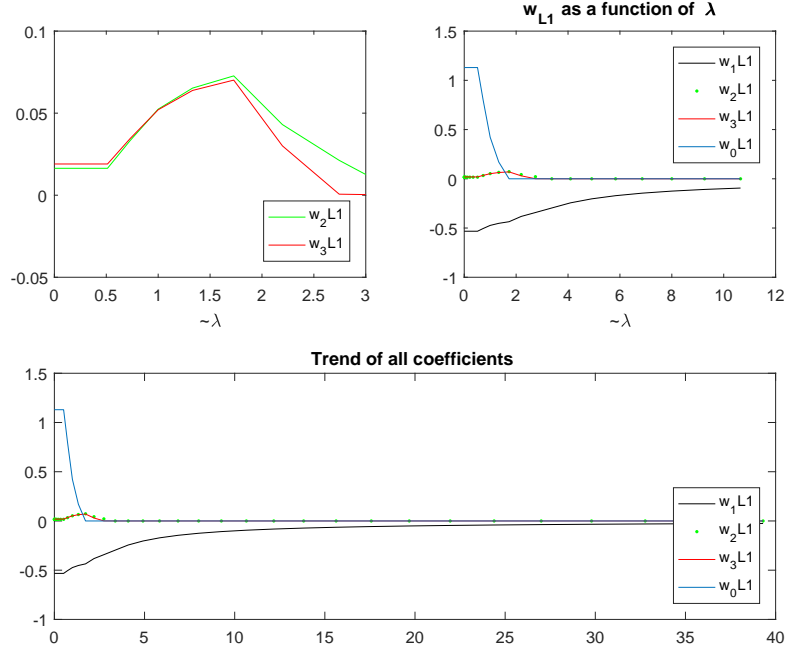


Figure 3: The RMS training and testing error for a wide range of $\frac{1}{\eta} \approx \lambda$.

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:

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

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

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

For the value of a given feature n , replaced by its class-conditional mean, $\mu_{c,n}$, this expression becomes 0; we further develop this by analyzing the log-odds:

$$\log \frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} = \log \frac{P(y = 1)}{P(y = 0)} + \log \frac{\frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1)}}{\frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x} - \mu_0)^T \Sigma^{-1}(\mathbf{x} - \mu_0)}} \quad (7)$$

Since the matrix sigma is shared, we can write:

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

Then we expand the exponent:

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

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 chosen $x_n = \mu_{c,n}$ 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

3.a how many parameters are there before and after

The initial model has five parameters, $\Theta_1, \Theta_{1,1}, \Theta_{0,1}, \Theta_{1,2}$ and $\Theta_{0,2}$. After duplication, we could think that we have seven parameters, the previous five and $\Theta_{1,3}, \Theta_{0,3}$, however since our third feature is a duplicate of the second, $\Theta_{1,3}, \Theta_{0,3}$ will behave exactly like $\Theta_{1,2}, \Theta_{0,2}$ and thus we still only have five parameters.

3.b What effect does this have on the decision boundary?

The decision boundary is given by:

$$w_0 + \sum_{i=1}^3 w_{i0} + \sum_{i=1}^3 (w_{i1} - w_{i0})x_i$$

where

$$\begin{aligned} w_0 &= \log \left(\frac{\Theta_1}{\Theta_0} \right) \\ w_{i1} &= \log \left(\frac{\Theta_{i1}}{\Theta_{i0}} \right) \\ w_{i0} &= \log \left(\frac{1 - \Theta_{i1}}{1 - \Theta_{i0}} \right) \\ \Theta_c &= P(y = c) \\ \Theta_{ic} &= P(x_i = c | y = c) \end{aligned} \quad (10)$$

Therefore and since feature 3 is a duplicate of feature 2, we could re-write the resulting boundary as:

$$w_0 + w_{10} + 2w_{20} + (w_{11} - w_{10})x_1 + 2(w_{21} - w_{20})x_2 \quad (11)$$

The preceding expression is similar to the boundary for the two-feature case but with a modified slope and bias, due to the weighting imposed on feature 2.

The greatest discrepancies for classification between the original model and the duplicate-feature model will occur when $\Theta_{20} \gg \Theta_{21}$ or $\Theta_{20} \ll \Theta_{21}$.

Additionally, since the inputs are binary, observations may only fall in one of four points on the plane³. Therefore, the worst effects of modifying the slope and bias can be felt when the boundary passes close to any one of these corners⁴.

Therefore, and because the sample space has strict constraints, we can say that Naive Bayes is fairly robust.

Question 4: Prove that for a gamma-normal prior and a normal distribution, the posterior will be gamma-normal, what are the parameters?

4.a Do so for the univariate case

Our prior distribution is given by:

$$\begin{aligned} \Pi(\mu, \lambda | \mu_0, \beta, a, b) &\equiv \Pi_1(\mu | \mu_0, (\beta\lambda)^{-1}) \Pi_2(\lambda | a, b) \\ &= (N)(\mu | \mu_0, (\beta\lambda)^{-1}) \text{Gam}(\lambda | a, b) \\ &= \left(\frac{\beta\lambda}{2\pi}\right)^{\frac{1}{2}} \exp\left\{-\frac{\beta\lambda}{2}(x_n - \mu)^2\right\} \frac{1}{\Gamma(a)} \lambda^{a-1} b^a \exp\{-b\lambda\} \end{aligned} \quad (12)$$

We know that our observations have a normal distribution consisting of N i.i.d. observations, so the likelihood $P(X|\mu, \lambda^{-1})$ can be expressed as a product:

$$\prod_{n=1}^N \left(\frac{\lambda}{2\pi}\right)^{\frac{1}{2}} \exp\left\{-\frac{\lambda}{2}(x_n - \mu)^2\right\} \quad (13)$$

The posterior is a product of the likelihood and the prior distribution. Thus the posterior can be expressed as:

$$\begin{aligned} \prod_{n=1}^N \left[\left(\frac{\lambda}{2\pi}\right)^{\frac{1}{2}} \exp\left\{-\frac{\lambda}{2}(x_n - \mu)^2\right\} \right] &\left(\frac{\beta\lambda}{2\pi}\right)^{\frac{1}{2}} \exp\left\{-\frac{\beta\lambda}{2}(\mu - \mu_0)^2\right\} \\ &\cdot \frac{1}{\Gamma(a)} \lambda^{a-1} b^a \exp\{-b\lambda\} \end{aligned}$$

³[0,0],[0,1],[1,0] and [1,1]

⁴e.g. a boundary with a slope slightly less than one with a negligible bias term would misclassify examples falling in the corner [1,1].

We absorb constant values not depending on μ , λ or β into a constant k and the previous product can equivalently be written as:

$$k \prod_{n=1}^N \left[\lambda^{\frac{1}{2}} \exp \left\{ -\frac{\lambda}{2} (x_n - \mu)^2 \right\} \right] \beta \lambda^{\frac{1}{2}} \exp \left\{ -\frac{\beta \lambda}{2} (\mu - \mu_0)^2 \right\} \lambda^{a-1} \exp \{-b\lambda\}$$

We further develop this:

$$\begin{aligned} & k \lambda^{a_0-1+\frac{N+1}{2}} \beta^{\frac{1}{2}} \exp \left\{ -\frac{\lambda}{2} \sum_{n=1}^N [(x_n - \mu)^2] \right\} \exp \left\{ -\frac{\beta \lambda}{2} (\mu - \mu_0)^2 \right\} \exp \{-b\lambda\} \\ &= k \lambda^{a_0-1+\frac{N+1}{2}} \beta^{\frac{1}{2}} \exp \left\{ -\frac{\lambda}{2} \sum_{n=1}^N [(x_n - \mu)^2] - \frac{\beta \lambda}{2} (\mu - \mu_0)^2 - b\lambda \right\} \end{aligned}$$

Since we can expand:

$$\begin{aligned} & \sum_{n=1}^N (x_n - \mu)^2 \\ &= \sum_{n=1}^N (x_n - \bar{x} + \bar{x} - \mu)^2 \\ &= \sum_{n=1}^N ((x_n - \bar{x}) - (\mu - \bar{x}))^2 \\ &= \sum_{n=1}^N (x_n - \bar{x})^2 - 2(x_n - \bar{x})(\mu - \bar{x}) + (\mu - \bar{x})^2 \\ & \quad \text{and since } \sum_{n=1}^N (x_n - \bar{x}) = 0 \\ &= \sum_{n=1}^N (x_n - \bar{x})^2 + (\mu - \bar{x})^2 \end{aligned}$$

Then we can arrive at:

$$k \lambda^{a_0-1+\frac{N+1}{2}} \beta^{\frac{1}{2}} \exp \left\{ -\frac{\lambda}{2} \left[\sum_{n=1}^N [(x_n - \bar{x})^2] + N(\mu - \bar{x})^2 + \beta(\mu - \mu_0)^2 \right] - b\lambda \right\} \quad (14)$$

Firstly, we notice that:

$$\sum_{n=1}^N [(x_n - \bar{x})^2] = \frac{N}{\lambda_d} \quad (15)$$

Where λ_d is the precision of the distribution that we have been given. We now focus on the terms of the exponent which are dependent on μ :

$$\begin{aligned}
& N(\mu + \bar{x})^2 + \beta(\mu - \mu_0)^2 \\
&= N\mu^2 + 2N\mu\bar{x} + N\bar{x}^2 + \beta\mu^2 - 2\beta\mu\mu_0 + \beta\mu_0^2 \\
&= (\beta + N)\mu^2 - 2\mu(\beta\mu_0 + N\bar{x}) + N\bar{x}^2 + \beta\mu_0^2 \\
&= (\beta + N) \left(\mu^2 - 2\mu \frac{(\beta\mu_0 + N\bar{x})}{\beta + N} + \left(\frac{\beta\mu_0 + N\bar{x}}{\beta + N} \right)^2 \right) + N\bar{x}^2 + \beta\mu_0^2 - \frac{(\beta\mu_0 + N\bar{x})^2}{\beta + N} \\
&= (\beta + N) \left(\mu - \left(\frac{\beta\mu_0 + N\bar{x}}{\beta + N} \right) \right)^2 + N\bar{x}^2 + \beta\mu_0^2 - \frac{(\beta\mu_0 + N\bar{x})^2}{\beta + N}
\end{aligned} \tag{16}$$

We turn our attention to the last three terms:

$$\begin{aligned}
& N\bar{x}^2 + \beta\mu_0^2 - \frac{(\beta\mu_0 + N\bar{x})^2}{\beta + N} \\
&= \frac{\beta N\bar{x}^2 + 2\beta N\bar{x}\mu_0 + \beta N\mu_0^2 - (\beta\mu_0 + N\bar{x})^2}{\beta + N} \\
&= \beta N \frac{(\bar{x} - \mu_0)^2}{\beta + N}
\end{aligned}$$

We plug this back into equation 16 and in turn insert equation 16 into equation 14:

$$k\lambda^{a_0-1+\frac{N+1}{2}}\beta^{\frac{1}{2}}\exp\left\{-\frac{\lambda}{2}\left(\frac{N}{\lambda_d} + (\beta + N)\left(\mu - \left(\frac{\beta\mu_0 + N\bar{x}}{\beta + N}\right)\right)^2 + \beta N \frac{(\bar{x} - \mu_0)^2}{\beta + N}\right) - b\lambda\right\}$$

We can rearrange the previous equation and modify k to be $k/\sqrt{1 + \frac{N}{\beta}}$:

$$\begin{aligned}
& \lambda^{a+\frac{N}{2}-1}\exp\left\{\left(-b + \frac{1}{2}\left(\frac{N}{\lambda_d} + \frac{\beta N(\bar{x} - \mu_0)^2}{\beta + N}\right)\right)\lambda\right\} \\
& \cdot ((\beta + N)\lambda)^{\frac{1}{2}}\exp\left\{\frac{\beta + N}{2}\lambda\left(\mu - \frac{\beta\mu_0 + N\bar{x}}{\beta + N}\right)^2\right\}k'
\end{aligned}$$

From where we can identify:

$$\begin{aligned}
\mu'_0 &= \left(\frac{\beta\mu_0 + N\bar{x}}{\beta + N}\right) \\
\beta' &= \beta + N \\
a' &= a + \frac{N}{2} \\
b' &= b + \frac{1}{2}\left(\frac{N}{\lambda_d} + \beta N \frac{(\bar{x} - \mu_0)^2}{\beta + N}\right)
\end{aligned}$$

Question 5: Implementation of Logistic Regression and Gaussian Naive-Bayes

A little bit about our data

Before plunging into the implementation task at hand, we take a look at some of the defining characteristics of our data set. With this heuristic purpose in mind, a routine was created to visualize some characteristics of our dataset. The code is included in the `plot_dyn.m` script.

```
for i = 1:32
    figure(3)
    plot(X(y==1, i), zeros(sum(y==1),1), 'g.')
    hold on
    plot(X(y==0, i), zeros(sum(y==0),1), 'r.')
    exes = min(X(:, i)) - 1:0.1:max(X(:, i)) + 1;
    plot(exes, normpdf(exes, mean(X(y == 1, i)),
        std(X(y == 1, i))), 'g')
    plot(exes, normpdf(exes, mean(X(y == 0, i)),
        std(X(y == 0, i))), 'r')
    title(['feature_', num2str(i), '_class_distributions'])
    hold off
    for j = i+1:32
        figure(4);
        plot(X((y == 1), i), X((y == 1), j), 'g+');
        hold on;
        plot(X((y == 0), i), X((y == 0), j), 'r.');
        hold off;
        title(['feature_', num2str(i),
            '_against_feature_', num2str(j)]);
        pause(0.1)
    end
end
```

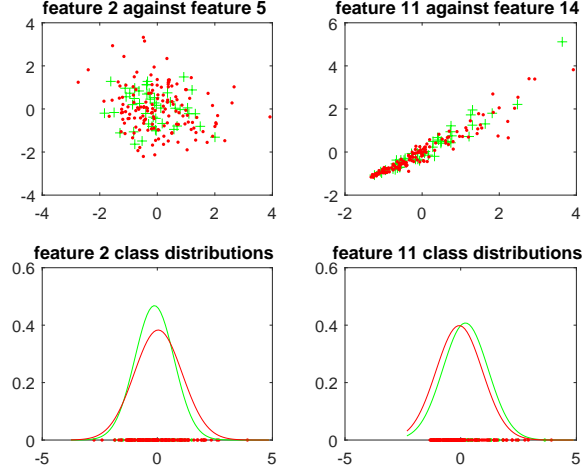


Figure 4: Some of the feature-pairs do not show a strong correlation (top-left), while other pairs (top-right) display a clear correlation. The class-distributions for single features are fairly close in most cases (bottom).

Firstly, we observe in figure 4 that the class-means and class variances of each one of our features are rather close, we have picked two examples of scatter-plots for feature pairs and two examples for the per-feature class-distribution. We can also observe that the data is nicely centered around 0 for all features and since the variances are not too different between all features⁵, then we can safely say that logistic regression would not benefit much from feature normalization.

A second characteristic that becomes quickly obvious is that some features are highly correlated, an example of this is plainly visible in figure 4 and we can also see this well represented in by the fact that the variance-covariance matrix of the inputs in figure 5 has large fairly values outside of the diagonal entries; this is not ideal for the naive assumption made by the gaussian-naive-bayes model.

Lastly, we note that there are 46 entries corresponding to class 1 out of 194 total entries. So if we were to arbitrarily classify all of our test data as negative, we would have a 76% rate of succesful classification.

⁵the avid reader is urged to run the included `plot_dyn.m` to cross-check this fact

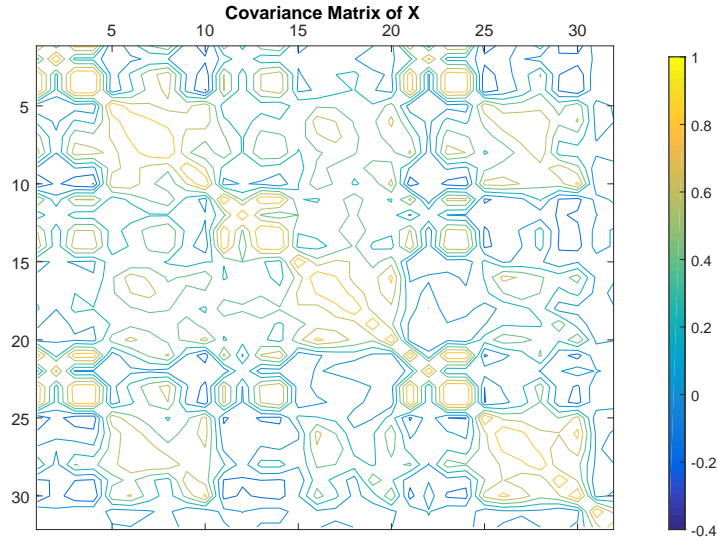


Figure 5: The variance-covariance matrix of the inputs shows large values outside of the diagonal entries.

5.a Logistic regression - gradient descent

For logistic regression, we chose to carry out a gradient descent implementation. The implementation is contained in the `LR_grad.m` function file with an auxiliary `find_alpha.m` function file.

Since gradient descent may get stuck on shallow local minima, iterating the learning process over a randomly-picked collection of 24 weight vectors was chosen. The aim of starting with several random initial weights being to ultimately choose the final vector that should yield the lesser error after all iterations of gradient descent should be done.

```

LR_grad.m:
function [ w, w_inits ] =

    % default to 24 random starting points
    if nargin < 4 || isempty(num_rand_inits)
        num_rand_inits = 24;
    end

    % if no random starting points are given, make some
    if nargin < 3 || isempty(w_inits)
        w_inits = rand(size(X,2), num_rand_inits);
    else
        num_rand_inits = size(w_inits(2));
    end

    w = w_inits;
    for i = 1: num_rand_inits
        done = false;
        prev_epsilon = realmax * ones(length(y),1);
        iterations = 1;
        % find a good learning rate for each vector
        lr = find_alpha(X, y, w(:,i));
        % Don't allow GD to run after it's found minima
        % nor for too long
        while ~done && (iterations < 2000)
            epsilon = y - (1./(1+exp(-X*w(:,i))));
            % if our prediction error is very small
            % or it's not really changing, we're done
            if (abs(sum(epsilon)) < .1) ...
                || abs(sum(epsilon)) ...
                    - sum(prev_epsilon)) < 0.01
                done = true;
            else % keep going down
                w(:,i) = w(:,i) + lr * (X'*epsilon);
                prev_epsilon = epsilon;
            end
            iterations = iterations + 1;
        end
    end
end

```

For each initial vector, an optimal learning rate is chosen from $\{\frac{1}{1}, \frac{1}{2}, \dots, \frac{1}{10}\}$. The learning rates remained fixed since implementing the Robbins-Monro conditions lead to a substantial degradation in performance - speedwise. The rationale to support this choice being that as the gradient approaches 0, the update to the \mathbf{w} vector will inherently be smaller.

The method for finding an optimal learning rate is implemented in the function `find_alpha.m`:

```
function alpha = find_alpha(X, y, w)
    lrs = 1./(1:10);
    errors = zeros(1, length(lrs));
    for i = 1: length(lrs)
        iterations = 1;
        lr = lrs(i);
        while iterations < 100
            epsilon = y - (1./(1+exp(-X*w)));
            w = w + lr * (X'*epsilon);
            iterations = iterations + 1;
        end
        errors(i) = sum(epsilon);
    end
    [dummy, idx] = min(abs(errors));
    alpha = lrs(idx);
end
```

In order to perform 10-fold cross-validation, we used matlab's `cvpartition` function, which yields homogenous grouping of entries per fold⁶. Matlab's `randperm` was also used to chose the entries randomly but without repetition for each fold. The following code from `A2.q5_driver.m` conatins the calls to `LR_grad.m`⁷.

```
...
folds_info = cvpartition(length(y), 'kfold', num_folds);
folds_idx = randperm(length(y));

num_rand_inits = 24;
features = 1:33;

for fold = 1:num_folds
    % just indexing for folds
    idxs_prev = 1:sum(folds_info.TestSize(1:(fold-1)));
    if ~isempty(idxs_prev)
        offset = idxs_prev(end);
    else
        offset = 0;
    end
    idxs_xcl = (1:folds_info.TestSize(fold))+offset;
    idx_after_skip = length(y) -
        (sum(folds_info.TestSize((fold + 1):end)) - 1);
```

⁶avoiding the final instance of a fold with 4 entries.

⁷The full version also contains calls to the naive Bayes Gaussian classifier, which will appear later on in section 5.b

```

idxs_next = idx_after_skip:length(y);
X_train = X(folds_idx([idxs_prev, idxs_next]),
            features);
X_test = X(folds_idx(idxs_xcl), features);
y_train = y(folds_idx([idxs_prev, idxs_next]));
y_test = y(folds_idx(idxs_xcl));

% Here's the actual call to logistic regression
[w, w_init] = LR_grad(X_train, y_train, [],
                    num_rand_inits);
w_grad(fold, :, :) = w;
for i=1:num_rand_inits
    errors_grad(fold, i) =
        sum(y_test ~= round(1
        ./ (1 + exp(- X_test*w(:,i))))) );
end
end
...

```

The choice of error was made so as to allow a proper comparison with the errors given by the Gaussian naive Bayes, or GNB hereafter. Since the implementation of GNB that was made for this assignment computes log-odds, the translation to a cross-entropy type error would have proven to be unwieldy. Thus the cost function was chosen to be of the form:

$$J_w(x_i) = \begin{cases} 0 & \text{if } \hat{y}_i = y_i \\ 1 & \text{if } \hat{y}_i \neq y_i \end{cases}$$

This is easily done via the following matlab command⁸:

```
min_err = min(abs(mean(errors_grad, 1)))}
```

⁸Note that only the minimum error of the 24 logistic regressions performed is taken into account for computing the error since this is one of the main reasons for having several random starting vectors.

5.b Implementation of GNB classifier

In order to setup our GNB algorithm, we separated the necessary actions into `gnb_train.m` and `gnb_predict.m`

The contents of `gnb_train.m`:

```
function [theta, mu_1, mu_0, Sigma] = gnb_train(X, y)
    theta = sum(y == 1)/length(y);
    idx_1 = y == 1;
    idx_0 = y == 0;
    mu_1 = mean(X(idx_1, :));
    mu_0 = mean(X(idx_0, :));
    x_to_mu = zeros(size(X));
    x_to_mu(idx_1, :) = X(idx_1, :)
        - repmat(mu_1, sum(idx_1), 1);
    x_to_mu(idx_0, :) = X(idx_0, :)
        - repmat(mu_1, sum(idx_0), 1);
    Sigma = (x_to_mu' * x_to_mu)./length(y);
end
```

We note that the estimated covariance matrix is shared between classes and can be estimated with the last three commands of the function, as per the class notes, lecture 3, slide 31. This effectively merges the covariance matrices of both classes.

The contents of `gnb_predict.m`:

```
function [log_odds, p_x_y_1, p_x_y_0] = ...
    gnb_predict(X, theta, mu_1, mu_0, Sigma)
    m = size(X, 1);
    n = size(X, 2);
    Sigma_inv = pinv(Sigma);
    norm_term = 1/((2*pi)^(n/2)*sqrt(det(Sigma)));
    X_to_mu_1 = (X - repmat(mu_1, m, 1));
    X_to_mu_0 = (X - repmat(mu_0, m, 1));

    alpha_1 =
        ((X_to_mu_1 * Sigma_inv * X_to_mu_1')
         .* eye(m)) * ones(m, 1);
    alpha_0 =
        ((X_to_mu_0 * Sigma_inv * X_to_mu_0')
         .* eye(m)) * ones(m, 1);
    p_x_y_1 = norm_term*exp(-alpha_1./2);
    p_x_y_0 = norm_term*exp(-alpha_0./2);
    log_odds = (alpha_0./2 - alpha_1./2) + ...
        repmat(log(theta/(1 - theta)),
            length(alpha_1), 1);
end
```

We call the attention to the reader to the use of `alpha_1 =* eye(m)`
`* ones(m,1);`. We prefer this to a per-observation loop since we are using a
matrix-operation-optimized software, this yields an equivalent result to per-
forming a prediction for each observation i , within a for loop; in the later case,
 $(\mathbf{x}_i - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_c)$ the result is a scalar. However, each i^{th} scalar from the
previous expression is contained in the i^{th} diagonal entry of the matrix obtained
from $(\mathbf{X} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}_c)$.

In order to run a comparison of logistic regression vs GNB over a single fea-
ture, suffice it to change the line: `features = 1:33;` to `features = [1,33];`.
We then add the following code within the `folds` loop:

```
...
% remove the column of ones, GNB needs no bias as it's
% already centered around the mean...
X_train(:,end) = [];
X_test(:,end) = [];
[theta, mu_1, mu_0, Sigma] = gnb_train(X_train, y_train);
[log_odds, p1, p0] =
    gnb_predict(X_test, theta, mu_1, mu_0, Sigma);
...
```

For logistic regression with gradient descent, using only the first feature of the
data-set, the average misclassification rate is about 4.8 per fold, depending
on the random permutations from which the folds were constructed. GNB
performed only slightly better, typically yielding about 4.6 misclassifications
per fold.

We recall from the end of the section regarding our data, that if we chose
to always classify as class 0, we should err 46 times out of 194. If our classifier
misclassified 4.6 samples per-fold during 10-fold classification it's erring at the
same rate as always choosing class 0. Still, however, GNB slightly outperforms
our logistic regression systematically.

5.c Compare Logistic regression and GNB over the whole feature-set

In order to classify over the whole feature-set, `features = [1,33];` needs to
be changed back to `features = 1:33;`. The accuracy of logistic regression fell
slightly to an average rate of misclassification of around 5.6 per-fold while GNB
continued to misclassify about 4.6 times per-fold on average, systematically
outperforming logistic regression.

The error rate remained on par with arbitrarily always choosing class 0,
however GNB actually did choose class 1 sometimes as we can see in figure 6. It
did, however, end up choosing class 0 fairly often as one would expect given the
 $\log \frac{P(y=1)}{P(y=0)}$ term of the log-odds equation⁹ and the proximity of the distributions
per class accross features.

⁹Which in this case is roughly -0.5

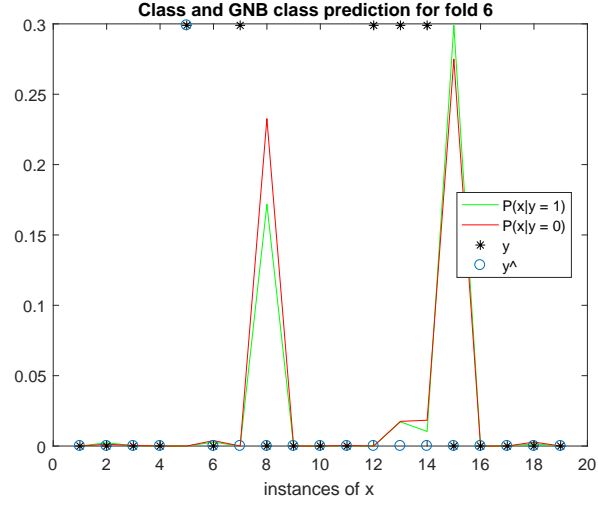


Figure 6: Data and predictions per instance. There is clearly a scaling issue with $P(x|y = c)$ for $c \in 0, 1$, however all we need is the ratio between $P(x|y = 1)$ and $P(x|y = 0)$. Note that instance observation 5 has been correctly classified as class 1.