

Homework 3

Subhankar Ghosh

Question 1

(a)

We will fit a *Linear*, *Polynomial* and a *Radial* kernel SVM on our digits data from ElemStatLearn package.

Five fold cross validation has been used to pick the best model

- Linear Kernel : $c = 0.0078125$ gives the best training accuracy of 99.07%
- Polynomial Kernel : The model with degree = 2, scale = 0.01 and $c = 0.5$ gives the best training accuracy of 99.53%
- Radial kernel : The model with sigma = 0.125 and $c = 2$ gives the best training accuracy of 79.09%

(b)

Table 1: Performance of different kernel SVM

Model_Name	Test_Accuracy
Linear SVM	0.9529412
Poly SVM	0.9500000
Radial SVM	0.1666667

(c)

Advantages of kernels used: - *Linear*: Faster to train than other kernels. Performs especially well when the separation boundary is linear. Works very well with multiple features. Performs very well when number of features is large. - *Polynomial*: Works very well when some polynomial space of the features have a linear separation boundary. Allows learning of non linear data models. Helps take care of interaction features as well. - *Radial*: Performs very well with non linear data, and when number of features is not very high. It can project features space into a higher dimensional feature space.

Disadvantages of kernels used: - *Linear*: Cannot be used if the decision boundary is non linear which is mostly the case in real life data. - *Polynomial*: Can be very slow as the degree of the polynomial increases. Can suffer from the problem of numerical instability. - *Radial*: Can be very slow and does not perform very well if the number of features is high.

Is cross-validation effective on this analysis? Is there any potential pitfall that you could think of?

Cross validation is effective in tuning the best parameter. But one potential pitfall can be the amount of time it takes to build each model and validate each of the many models it has to train.

Question 2

(a) My own linear discriminate analysis (LDA) code gave an accuracy of **0.82** on the test data.

The LDA from the MASS package also gave an accuracy of **0.82** on the test data.

(b) My own linear discriminate analysis (LDA) code gave an accuracy of **0.835** on the test data with an $\alpha = 0.05$.

Dual form:
$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j x_i^T x_j$$

subject to: $\alpha_i \geq 0 \quad i=1, \dots, n$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

This can be written as
$$\min_{\alpha_i} - \sum_{i=1}^n \alpha_i + \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j x_i^T x_j$$

conditions remain same.

Solve qp solves the problem
$$\min \frac{1}{2} \beta^T D \beta - d^T \beta$$

subject to: $A^T \beta \geq b_0$

$$D = (Y^* Y^T) (X^* X^T)$$

$$d = [1, 1, \dots, 1]^T_n$$

$$A^T = \begin{bmatrix} y_1 & y_2 & \dots & y_n \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ 0 & 0 & 1 & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \quad b_0 = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$(n+1) \times n \qquad n \times 1 \times 1$

Figure 1: Formulation of Dual form

(c) Regularized QDA has an advantage over LDA because in this case we do not need to study the structure of the data in order to decide if we should apply LDA or QDA. It automatically decides the appropriate value of α thus deciding whether a reduction to a particular form (LDA or QDA) is required or not.

QDA is unlikely to produce satisfactory results unless the ratio of the class sample sizes is large relative to the number of variables posed in the problem, and it has no advantages compared to LDA except when the class covariance matrices are quite different. In situations where the class covariance matrices are similar, LDA is likely to give improved results compared to QDA, because a single covariance matrix can be employed in the classification rule, thus reducing the number of estimates which need to be computed. LDA should therefore be able to produce better results, when the sample sizes are smaller.

Regularized QDA always gives results equivalent to or better than LDA and QDA since it automatically decides the appropriate value of α thus deciding whether a reduction to a particular form (LDA or QDA) is required or not. The only downside to Regularized QDA is that it takes a lot of time.

Question 3

(a)

```
## [1] "Coefficients from my SVM"
```

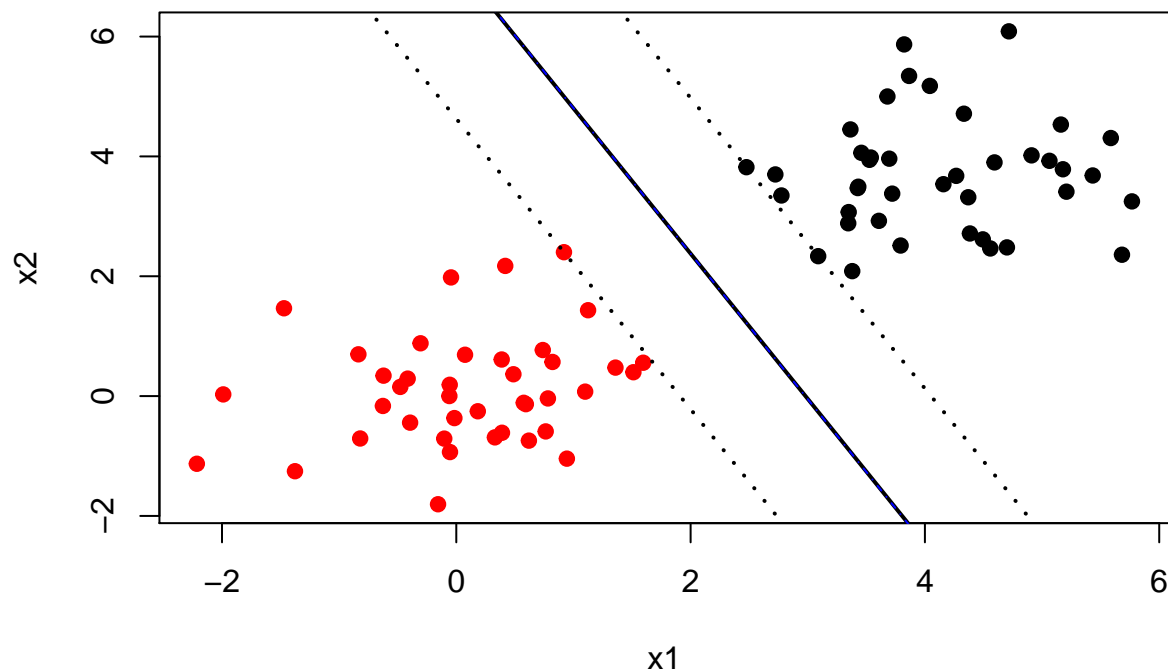


Figure 2: Comparison plot of my svm vs package. Blue:SVM Package, Black: My LDA

```
## [1] "Beta0 = 2.78237298573782  Beta values = -0.933425983836912, -0.384981990665094"
```

```
## [1] "Coefficients from package SVM"
```

```
## [1] "Beta0 = 2.78195312036252  Beta values = -0.933873966286481, -0.384495654907918"
```

(b)

Plot showing decision boundary we got from my SVM with that we get from the package

(c)

(d)

$$\text{Dual form: } \max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j x_i^T x_j$$

$$\text{s.t. } \alpha_i \geq 0 \quad i=1,2,\dots,n \quad \sum_{i=1}^n \alpha_i y_i = 0.$$

$$\alpha_i \leq c \quad i=1,2,\dots,n$$

Converting to minimization problem:

$$\min_{\alpha_i} - \sum \alpha_i + \frac{1}{2} \sum y_i y_j \alpha_i \alpha_j x_i^T x_j$$

$$0 \leq \alpha_i \leq c \quad i=1,2,\dots,n$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

$$D = (y y^T) (x x^T)$$

$$d = [1, 1, \dots, 1]_{1 \times n}$$

$$A^T = \begin{bmatrix} y_1 & y_2 & \dots & y_n \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ -1 & 0 & \dots & 0 \\ \vdots & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & -1 \end{bmatrix}_{(2n+1, n)}$$

$$b_0 = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ -c \\ -c \\ \vdots \\ -c \end{bmatrix}_{(2n+1) \times 1}$$

Figure 3: Formulation of dual form

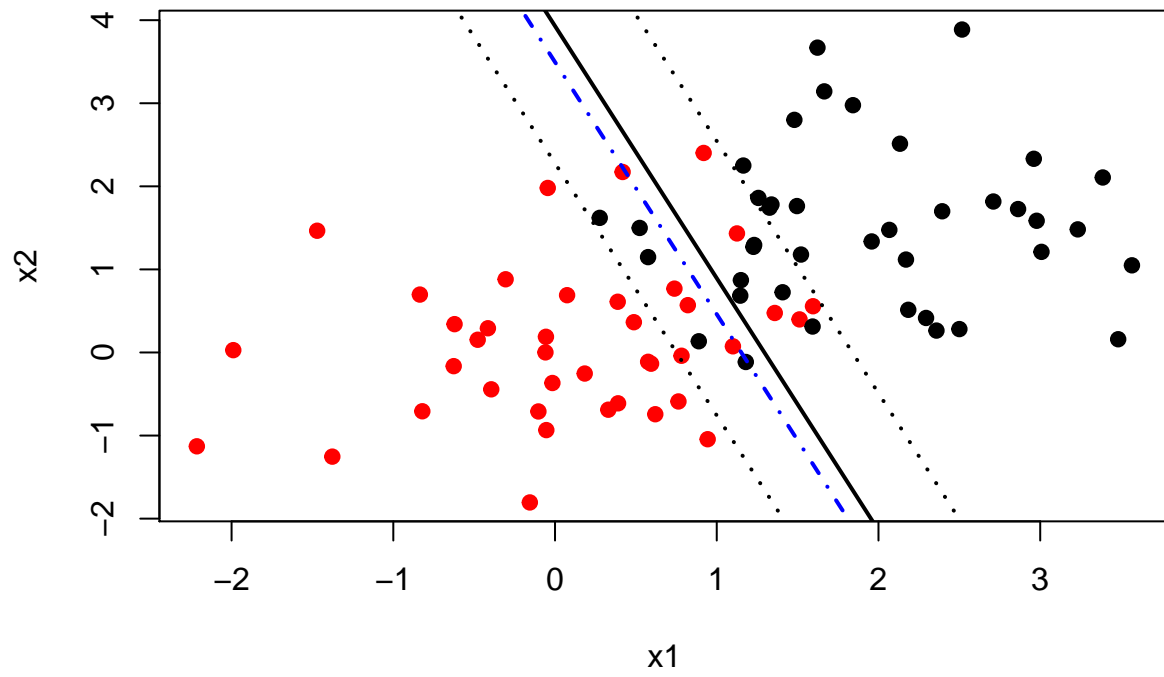


Figure 4: Comparison plot of my svm vs package. Blue:SVM Package, Black: My Regularized QDA

$$(3) (c) \quad \Phi(x) = (\phi_1(x), \phi_2(x), \dots, \phi_n(x))$$

Transforming the feature space from \mathcal{T} to m .

$$f(x) = \langle \Phi(x), \beta \rangle$$

$$F(x) \text{ is given by } f(x) = \beta_0 + x^T \beta$$

$$\text{Substituting } x \text{ by } \Phi(x) \Rightarrow f(x) = \beta_0 + \Phi(x)^T \beta, \quad \beta \in \mathbb{R}^m$$

$$\text{we get: } \beta - \sum_{i=1}^n y_i \alpha_i = 0 \quad \text{and} \quad \sum_{i=1}^n \alpha_i y_i = 0$$

Now we have the problem of

$$\min \|\beta\|^2 \text{ s.t. } y_i f(x_i) \geq 1$$

$$f(x) = \sum \alpha_i y_i \Phi(x)^T \Phi(x) + \beta_0$$

$K(x, z)$: kernel function.

$$K(x, z) = \langle \Phi(x), \Phi(z) \rangle$$

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i y_i \langle \Phi(x), \Phi(z) \rangle$$

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i y_i K(x, z)$$

Figure 5: Formulation of new decision function

(D) Dual form: $\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \langle \phi(x_i) \phi(x_j) \rangle$
s.t. $0 \leq \alpha_i \leq c$, $\sum_{i=1}^n \alpha_i y_i = 0$

Corresponding kernel

$\max_{\alpha_i} \sum \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j k(y_i, y_j)$ same constraint as above

\therefore Variables of solve. QP will be:-

$D = \begin{bmatrix} y_1 y_1 k(x_1, x_1) & \dots & y_1 y_n k(x_1, x_n) \\ \vdots & & \vdots \\ y_n y_1 k(x_n, x_1) & \dots & y_n y_n k(x_n, x_n) \end{bmatrix}_{n \times n}$ $D_{ij} = y_i y_j k(x_i, x_j)$
 $d_{n \times 1} = [1 \dots 1]^T_{1 \times n}$

$A = \begin{bmatrix} y_1 & 1 & 0 & 0 & \dots & 0 & -1 & 0 & \dots & 0 \\ y_2 & 0 & 1 & & & & 0 & -1 & & \\ \vdots & & & \ddots & & & & & \ddots & \\ y_n & 0 & & & 1 & & 0 & & & -1 \end{bmatrix}_{n \times (R_n + 1)}$ $b_0 = [0 \ 0 \dots 0 \ -c \dots -c]$
 $\sum \alpha_i y_i = 0$

Figure 6: Formulation of dual problem