

Q1.1

$$\text{minimize}_{w,b} \lambda \|w\|^2 + \sum_{i=1}^n (w^T x_i + b - y_i)^2$$

$$\text{minimize } \lambda w^T w + (Y - X^T w)^2$$

$$L = \lambda w^T w + (Y - X^T w)^T (Y - X^T w)$$

$$\frac{\partial L}{\partial w} = 2\lambda w - 2X(Y - X^T w) = 0$$

$$\lambda w - XY + XX^T w = 0$$

$$\lambda w + XX^T w = XY$$

$$w(XX^T + \lambda I) = XY$$

$$w = (XX^T + \lambda I)^{-1} XY \quad \text{Hence Proved.}$$

Q1.2)

$$d_{(i)} = d - x_i y_i$$

$$C_{(i)} = C - x_i x_i^T$$

(Q1.3)

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^TA^{-1}}{1 + v^TA^{-1}u}$$

$$C = xx^T + \lambda I^T$$

$$C_{(i)}^{-1} = (C - x_i x_i^T)^{-1}$$

$$C_{(i)}^{-1} = C^{-1} + \frac{C^{-1}x_i x_i^T C^{-1}}{1 - x_i^T C^{-1} x_i}$$

Q1.4

$$w = C^{-1}d$$

$$w_i = C_i^{-1}d_i$$

$$C_i^{-1} = (C - x_i x_i^T)^{-1}$$

$$d_i = d - x_i y_i$$

$$w_i = \begin{bmatrix} C^{-1} + \frac{C^{-1}x_i x_i^T C^{-1}}{1 - x_i^T C^{-1} x_i} \end{bmatrix} [d - x_i y_i]$$

$$w_i = C^{-1}d - C^{-1}x_i y_i + \frac{C^{-1}x_i x_i^T C^{-1}d - C^{-1}x_i x_i^T C^{-1}x_i y_i}{1 - x_i^T C^{-1} x_i}$$

$$= w - \frac{C^{-1}x_i y_i + C^{-1}x_i y_i x_i^T C^{-1}x_i + C^{-1}x_i x_i^T C^{-1}d - C^{-1}x_i x_i^T C^{-1}x_i y_i}{1 - x_i^T C^{-1} x_i}$$

$$= w + \frac{C^{-1}x_i [-y_i + y_i x_i^T C^{-1}x_i + x_i^T C^{-1}d - y_i x_i^T C^{-1}x_i]}{1 - x_i^T C^{-1} x_i}$$

$$w_i = w + \frac{C^{-1}x_i [-y_i + x_i^T w]}{1 - x_i^T C^{-1} x_i}$$

Hence proved

Q1.5

From previous solution:

$$\bar{w}_i = \bar{w} + (C^{-1}\bar{x}_i) \frac{[-y_i + \bar{x}_i^T \bar{w}]}{1 - \bar{x}_i^T C^{-1} \bar{x}_i}$$

To prove:  $\bar{w}_i^T \bar{x}_i - \bar{y}_i = \frac{\bar{w}^T \bar{x}_i - \bar{y}_i}{1 - \bar{x}_i^T C^{-1} \bar{x}_i}$

$$\bar{w}_i^T \bar{x}_i - \bar{y}_i = \left[ \bar{w} + \frac{C^{-1}\bar{x}_i [-y_i + \bar{x}_i^T \bar{w}]}{1 - \bar{x}_i^T C^{-1} \bar{x}_i} \right]^T \bar{x}_i - \bar{y}_i$$

$$= \bar{w}^T \bar{x}_i - \frac{\bar{w}^T \bar{x}_i \bar{x}_i^T C^{-1} \bar{x}_i}{1 - \bar{x}_i^T C^{-1} \bar{x}_i} - \frac{y_i \bar{x}_i^T C^{-1} \bar{x}_i}{1 - \bar{x}_i^T C^{-1} \bar{x}_i} + \frac{\bar{x}_i^T \bar{w} C^{-1} \bar{x}_i}{1 - \bar{x}_i^T C^{-1} \bar{x}_i} - \bar{y}_i$$

$$\boxed{\bar{w}_i^T \bar{x}_i - \bar{y}_i = \frac{\bar{w}^T \bar{x}_i - \bar{y}_i}{1 - \bar{x}_i^T C^{-1} \bar{x}_i}}$$

Hence proved.



Q1.6

$$w_{(k)}^T x_i - y_i = \frac{w^T x_i - y_i}{1 - x_i^T C^{-1} x_i}$$

$$C = (xx^T + \lambda I)$$

where dimensions of  $x$  are  $(k+1) \times n$

$y$  are  $n \times 1$

$w$  are  $(k+1) \times 1$

$C$  are  $(k+1) \times (k+1)$

where  $k$  is the number of features &  $n$  is the number of instances

The complexities of computing

$$w^T x_i = O(k+1)$$

$$C = O(k+1)^2$$

$$C^{-1} = O(k+1)^3$$

$$x_i^T C^{-1} x_i = O(k+1)^2$$

$$\text{Overall Complexity} = O[(k+1) + (k+1)^2 n + (k+1)^3 + (k+1)^2]$$

This is for 1 instance, for  $n$  instances it will be

$$O[n(k+1) + n^2(k+1)^2 + n(k+1)^3 + n(k+1)^2] + O(n(k+1)^3 + n(k+1)^2)$$

as  $C$  &  $C^{-1}$  are calculated just once

$$\therefore \text{Overall Complexity} = O[(k+1)^3 + n(k+1)^2] \quad \text{--- (1)}$$

From eq(1) in the question paper, we have

$$\bar{w} = C^{-1}d$$
$$\bar{w} = (\bar{x}\bar{x}^T + \lambda I)^{-1} \bar{x}y$$

where dimensions of  $x$  are  $(k+1)(n-1)$   
 $y$  are  $(n-1) \times 1$   
 $w$  are  $(k+1) \times 1$

The complexity for computing  $\bar{x}\bar{x}^T$  is  $(n-1)(k+1)^2 \sim n(k+1)^2$   
inverse of  $\bar{x}\bar{x}^T + \lambda I$  is  $(k+1)^3$   
multiplying  $\bar{x}y$  is  $(k+1)(n-1) \sim n(k+1)$   
multiplying inverse  $\bar{x}y$  is  $(k+1)^2$

$$\text{Total complexity} = O(n(k+1)^2 + (k+1)^3 + n(k+1) + (k+1)^2)$$

$$\text{But for } n \ll OCV, \text{ time complexity} = O(n^2(k+1)^2 + n(k+1)^3 + n^2(k+1) + n(k+1)^2)$$

$$\text{Overall complexity} = O(n^2(k+1)^2 + n(k+1)^3) \quad \text{--- (ii)}$$

From eq (ii) & (i), it can be seen that algorithmic complexity of computing LOOCV error using the formula in question 1.5 is less than computing it the usual way.

$$Q2.1 \quad P(Y=y|X) = \underset{y}{\operatorname{argmax}} P(Y=y) \cdot P(X|Y=y)$$

$$= \underset{y}{\operatorname{argmax}} P(Y=y) \cdot \prod P(x_i|Y=y)$$

$$= \propto P(Y=y) \cdot \prod P(x_i|Y=y) \quad \text{--- (1)}$$

Since  $x_1$  &  $x_2$  are independent of each.

Since  $x_2$  is a continuous variable, we have

$$P(X=x_2|Y=y) = \frac{1}{\sigma_y \sqrt{2\pi}} \exp \left[ -\frac{(x_2 - \mu_y)^2}{2\sigma_y^2} \right]$$

But for  $y=0$  &  $y=1$ , we have

$$P(X=x_2|Y=0) = \frac{1}{\sigma_0 \sqrt{2\pi}} \exp \left[ -\frac{(x_2 - \mu_0)^2}{2\sigma_0^2} \right]$$

$$P(X=x_2|Y=1) = \frac{1}{\sigma_1 \sqrt{2\pi}} \exp \left[ -\frac{(x_2 - \mu_1)^2}{2\sigma_1^2} \right]$$

So, we have 4 parameters to estimate for  $P(X=x_2|Y=y)$  those are  $\sigma_0$ ,  $\sigma_1$ ,  $\mu_0$ ,  $\mu_1$ .

Now, since  $x=x_1$  is a boolean variable

$$P(X=x_1|Y=y) = \alpha_y^{x_1} (1-\alpha_y)^{(1-x_1)}$$

But for  $y=0$  &  $y=1$ , we have

$$P(X=x_1|Y=0) = \alpha_0^{x_1} (1-\alpha_0)^{(1-x_1)}$$

$$P(X=x_1|Y=1) = \alpha_1^{x_1} (1-\alpha_1)^{(1-x_1)}$$



So, for  $x=x_i$ , we have 2 parameters to estimate for  $P(x=x_i | Y=y)$ , the parameters are  $\theta_0$  &  $\theta_1$ .

$\therefore$  Eqn becomes

$$P(Y=y|x) = \alpha \cdot P(Y=y) \cdot [\theta_y^{x_i} (1-\theta_y)^{(1-x_i)}] \left[ \frac{1}{\sigma_y \sqrt{2\pi}} \exp \left[ -\frac{(x_i - \mu_y)^2}{2\sigma_y^2} \right] \right]$$

Q2.2

$$1) P(Y=1|X) = \frac{P(X|Y=1) \cdot P(Y=1)}{P(X|Y=1) \cdot P(Y=1) + P(X|Y=0) \cdot P(Y=0)}$$

$$= \frac{1}{1 + \frac{P(X|Y=0) \cdot P(Y=0)}{P(X|Y=1) \cdot P(Y=1)}}$$

$$= \frac{1}{1 + \exp\left(\ln\left(\frac{P(X|Y=0) \cdot P(Y=0)}{P(X|Y=1) \cdot P(Y=1)}\right)\right)}$$

Let  $P(Y=0) = \gamma$ , we get

$$= \frac{1}{1 + \exp\left[\ln \frac{\gamma}{1-\gamma} + \sum_i \ln \left(\frac{P(x_i|Y=0)}{P(x_i|Y=1)}\right)\right]}$$

As  $x_1, x_2, \dots, x_d$  is a vector of ~~Boolean~~ <sup>Boolean</sup> variables, they follow binomial distribution

$$P(x_i|Y=0) = \theta_{i0}^{x_i} (1-\theta_{i0})^{(1-x_i)}$$

$$P(x_i|Y=1) = \theta_{i1}^{x_i} (1-\theta_{i1})^{(1-x_i)}$$

$$\therefore P(Y=1|X) = \frac{1}{1 + \exp\left[\ln \frac{\gamma}{1-\gamma} + \sum_i \ln \frac{\theta_{i0}^{x_i} (1-\theta_{i0})^{(1-x_i)}}{\theta_{i1}^{x_i} (1-\theta_{i1})^{(1-x_i)}}\right]}$$

$$= \frac{1}{1 + \exp \left[ \ln \frac{\gamma}{1-\gamma} + \sum_i x_i \ln \frac{\theta_{i0}}{\theta_{i1}} + (1-x_i) \ln \frac{(1-\theta_{i0})}{(1-\theta_{i1})} \right]}$$

$$= \frac{1}{1 + \exp \left[ \ln \frac{\gamma}{1-\gamma} + \sum_i \ln \frac{(1-\theta_{i0})}{(1-\theta_{i1})} + \sum_i x_i \left[ \ln \frac{\theta_{i0}}{\theta_{i1}} - \ln \frac{(1-\theta_{i0})}{(1-\theta_{i1})} \right] \right]}$$

$$\text{Let } w_0 = \ln \frac{\gamma}{1-\gamma} + \sum_i \ln \frac{(1-\theta_{i0})}{(1-\theta_{i1})}$$

$$w_i = \ln \frac{\theta_{i0}}{\theta_{i1}} - \ln \frac{(1-\theta_{i0})}{(1-\theta_{i1})}$$

$$P(Y=1|x) = \frac{1}{1 + \exp(w_0 + \sum_i w_i x_i)}$$

which is similar to Logistic Regression.

Q3.1.1

$$\max_{\alpha} \sum_{j=1}^n \alpha_j - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i \alpha_i y_j \alpha_j x_i x_j$$

$$\min_{\alpha} - \sum_{j=1}^n \alpha_j + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i \alpha_i y_j \alpha_j x_i x_j$$

$$\text{s.t.} \sum_{j=1}^n y_j \alpha_j = 0$$

$$0 \leq \alpha_j \leq c \quad \forall j$$

$$\min_x \frac{1}{2} x^T H x + f^T x$$

$$\text{s.t.} Ax \leq b$$

$$A_{\text{eq}} \cdot x = b_{\text{eq}}$$

$$lb \leq x \leq ub$$

$$x = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix}_{n \times 1}$$

$$\text{W.K.T.} \sum_{j=1}^n y_j \alpha_j = 0$$

$$\text{Also } A_{\text{eq}} \cdot x = b_{\text{eq}}$$

so.

$$A_{\text{eq}} = [y_1 \dots y_n]_{1 \times n}$$

$$b_{\text{eq}} = 0$$

$$\text{Also } 0 \leq \alpha_j \leq c \quad \forall j$$

$$\text{So } lb \leq x \leq ub$$

$$lb = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}_{n \times 1}$$

$$ub = \begin{bmatrix} c \\ \vdots \\ c \end{bmatrix}_{n \times 1}$$

$$A = [ ]$$

$$b = [ ]$$

$$H = x_i x_j y_i y_j$$



Q3.2

$$L_i = \frac{1}{2n} \sum_{j=1}^k \|w_j\|^2 + c L(w, x_i, y_i)$$

$$L(w, x_i, y_i) = \max \{ w_{\hat{y}_i}^T x_i - w_{y_i}^T x_i + 1, 0 \}$$

$$\hat{y}_i = \operatorname{argmax}_{j \neq y_i} w_j^T x_i$$

1) Subgradient of  $L_i$  wrt  $w_{y_i}$

$$\text{If } w_{\hat{y}_i}^T x_i - w_{y_i}^T x_i + 1 \leq 0$$

$$\text{then } \frac{\partial L_i}{\partial w_{y_i}} = \frac{w_{y_i}}{n}$$

else

$$\frac{\partial L_i}{\partial w_{y_i}} = \frac{w_{y_i}}{n} - c x_i$$

2) Subgradient of  $L_i$  wrt  $w_{\hat{y}_i}$

$$\text{If } w_{\hat{y}_i}^T x_i - w_{y_i}^T x_i + 1 \leq 0$$

$$\text{then } \frac{\partial L_i}{\partial w_{\hat{y}_i}} = \frac{w_{\hat{y}_i}}{n}$$

$$\text{else } \frac{\partial L_i}{\partial w_{\hat{y}_i}} = \frac{w_{\hat{y}_i}}{n} + c x_i$$

3 Subgradient of  $L_i$  wrt  $w_j$  for  $j \neq x_i$  &  $j \neq \hat{y}_i$

$$\frac{\partial L}{\partial w_j} = \frac{w_j}{n}$$

### **Question 3.1:**

For  $C = 10$  and  $\epsilon = 0.1$ :

**Accuracy:** 0.978202

**Objective function value:** 112.1461

**Confusion matrix :**

179 4

4 180

**Number of support vectors:** 119

For  $C = 0.1$  and  $\epsilon = 0.0001$ :

**Accuracy:** 0.907357

**Objective function value:** 24.7648

**Confusion matrix :**

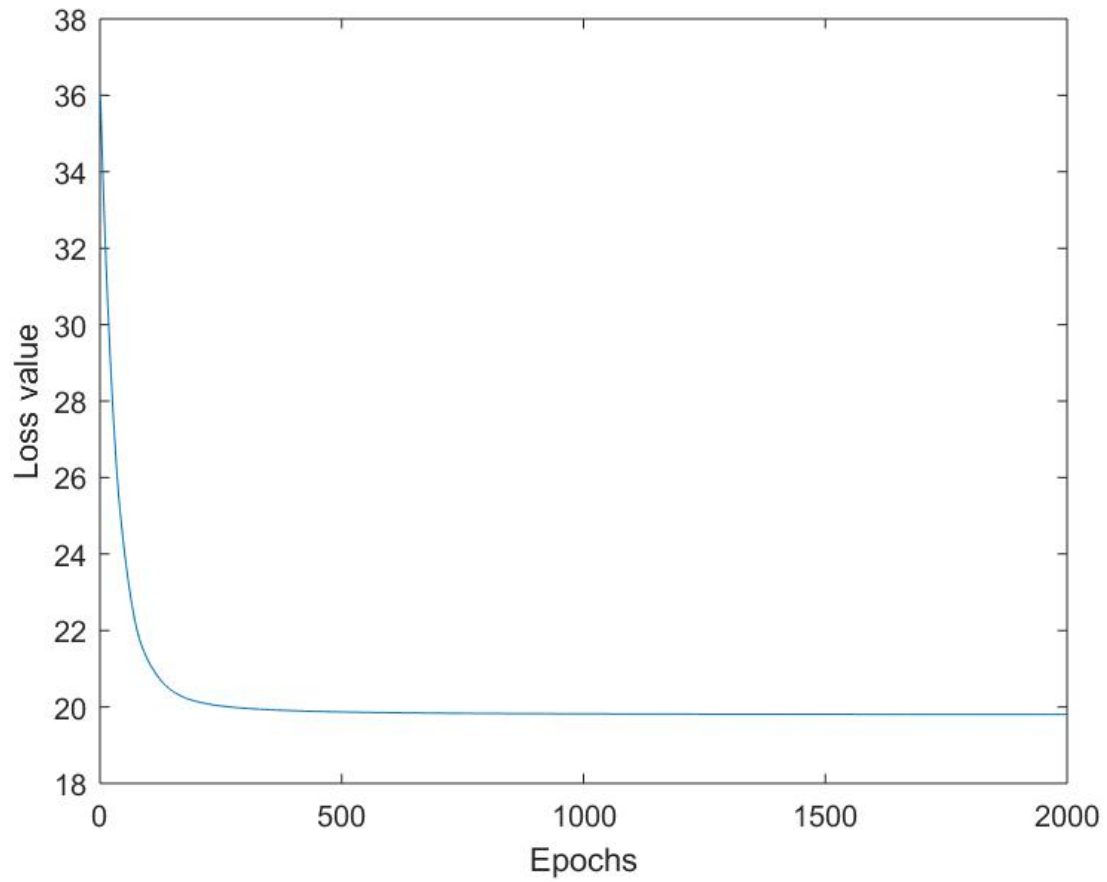
181 32

2 152

**Number of support vectors:** 339

### Question 3.2.5

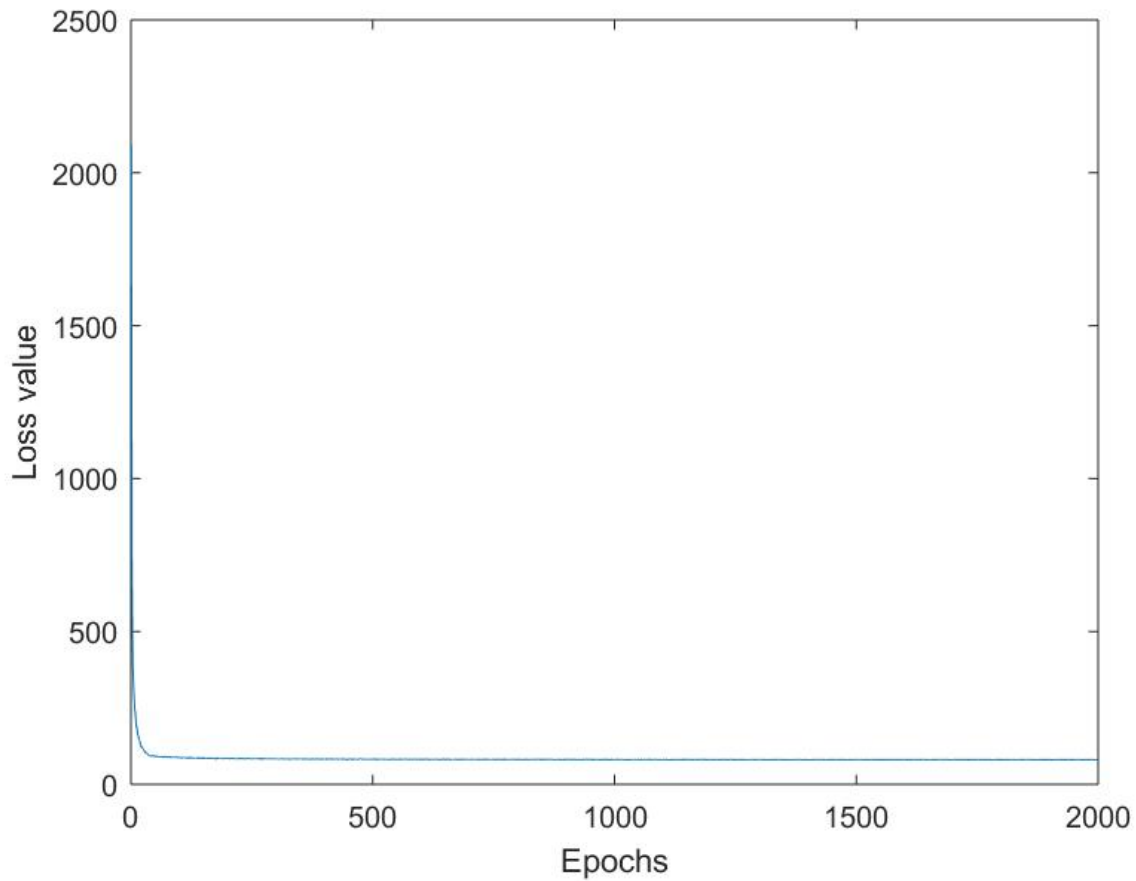
For  $c = 0.1$ ,  $\text{ita}_0 = 1$ ,  $\text{ita}_1 = 100$ ,  $\text{iterations} = 2000$ , the plot of loss values is



**Objective function value:** 19.8035



For  $C = 10$



**Objective function value: 80.7702**

The objective function value calculated using stochastic gradient descent is less than the objective function value calculated using quadratic programming.

### Question 3.2.6

For  $C = 0.1$

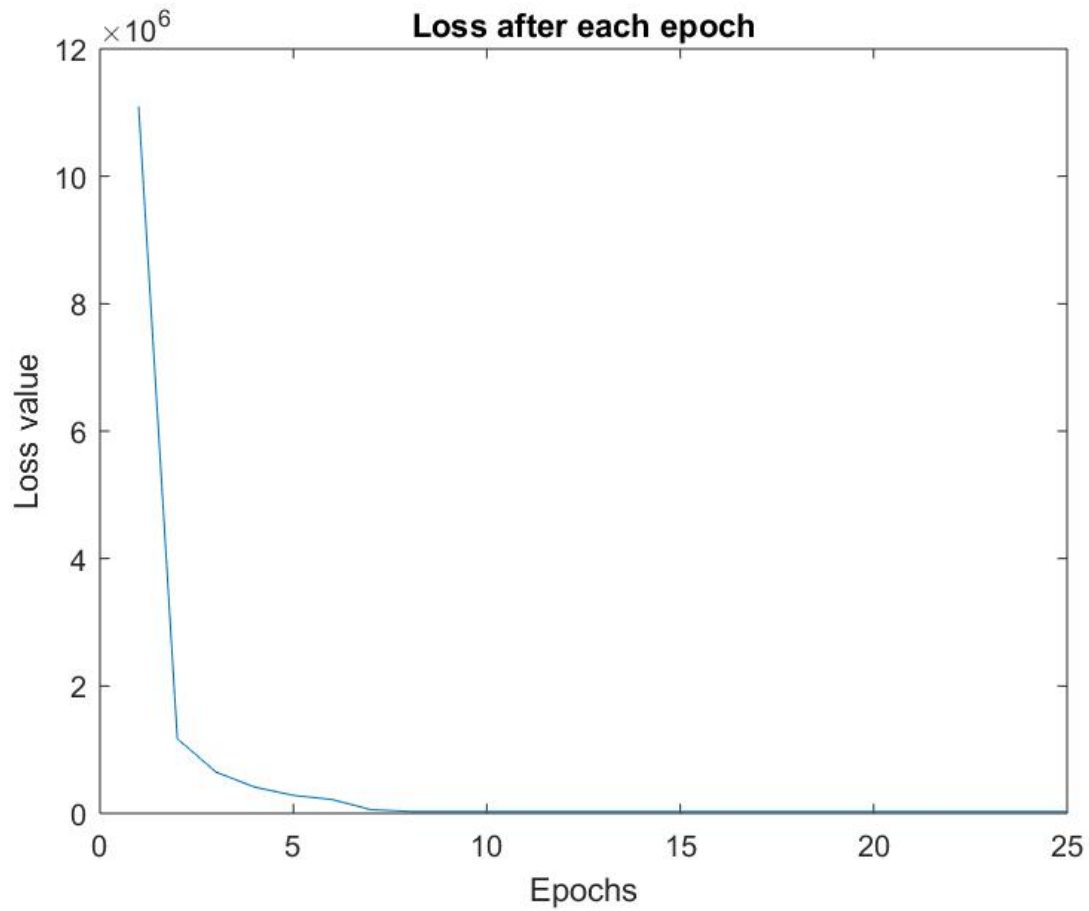
- a) **Prediction error on training data: 0.0359**
- b) **Prediction error on validation data (test error): 0.0572**
- c) **L2 norm of weights and its square: 16.1100**

For  $C = 10$

- a) **Prediction error on training data: 0.0055**
- b) **Prediction error on validation data (test error): 0.0354**

c) **L2 norm square of w:** 121.0098

### Question 3.2.7



The values on which I got the best result are:  $C = 8$ , epochs = 25, ita\_1 = 1000, ita\_0 = 1

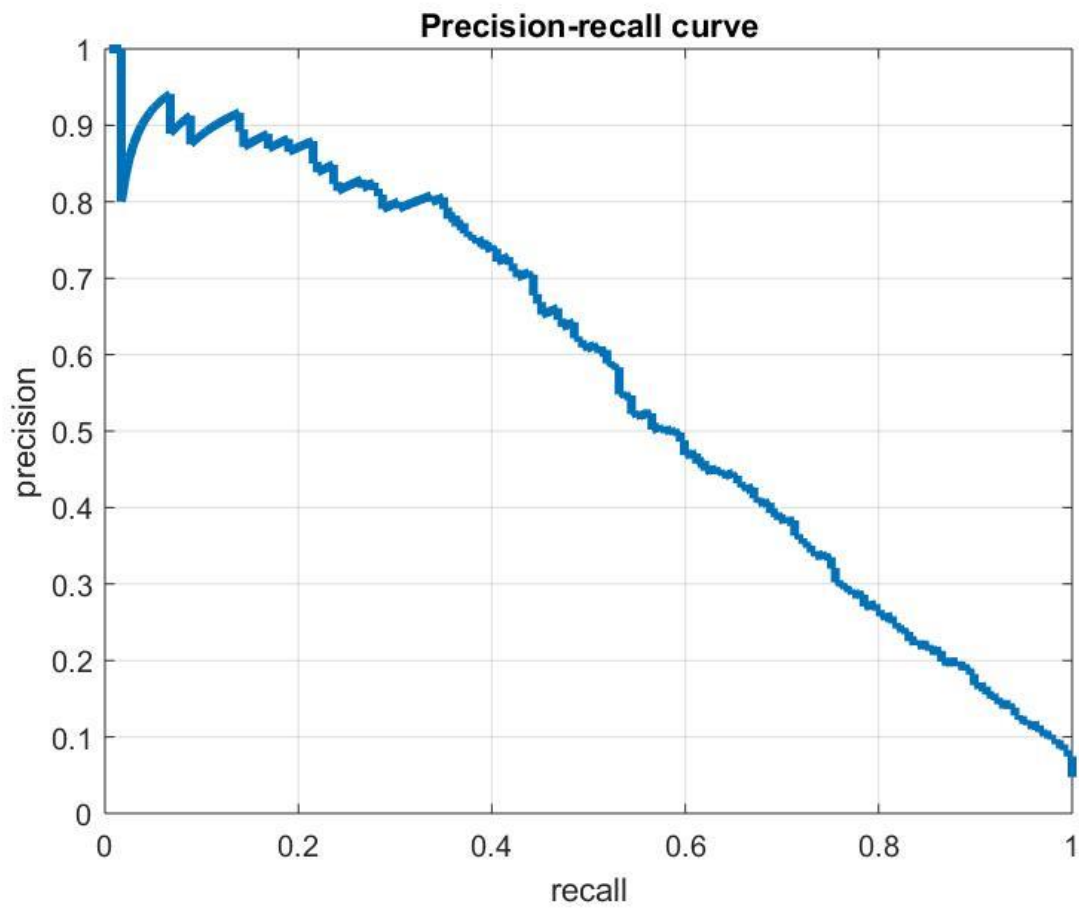
**Kaggle rank:** 8

**Kaggle score:** 0.79780

**Question 4.4.1 :**

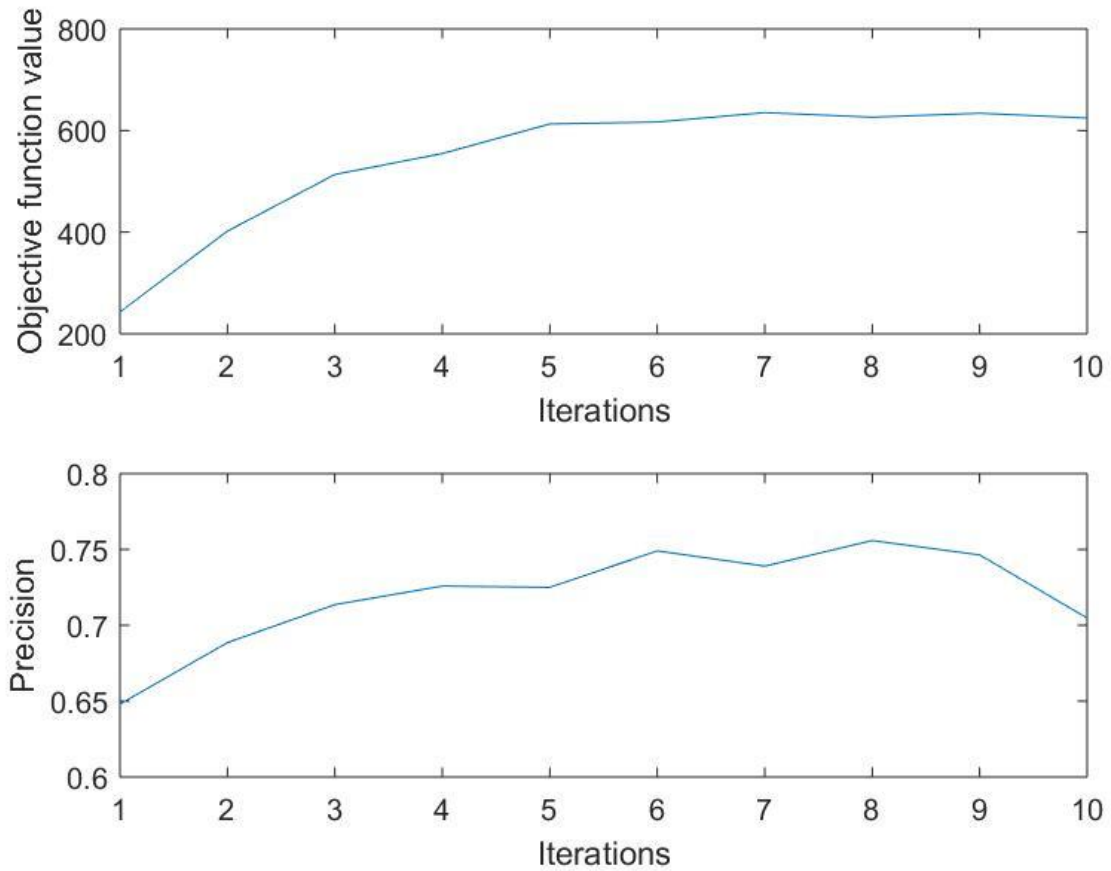
**AP values :** 0.578487068694247

**Precision recall curve:**



### Question 4.4.3

Plot of Objective values and the APs



### Question 4.4.4:

Value of AP: 79.63

(Best I achieved was 81.45)