

CSE 512 - Machine Learning HW 2

111462188 – Mihir Chakradeo

Q.1 Ridge Regression and LOOCV

QUESTION 1 - RIDGE REGRESSION AND LOOCV.

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

(1.1) $\bar{w} = [w; b]$, $\bar{x} = [x; 1]$, $\bar{I} = [I_K, 0_K; 0_K^T, 0]$, $C = \bar{x}\bar{x}^T + \lambda \bar{I}$
 $d = \bar{x}y$. Show that the solution of Ridge Regression is: $\bar{w} = C^{-1}d$.

Ans. $\min_{w, b} \lambda \|\bar{w}\|^2 + \sum_{i=1}^n (\bar{w}^T \bar{x}_i - y_i)^2.$

$$\Rightarrow \min_{w, b} \lambda \|\bar{w}\|^2 + \|\bar{x}^T \bar{w} - Y\|_2^2$$

Take gradient wrt \bar{w} and set to 0.

$$2\lambda \bar{w} + 2\bar{x}(\bar{x}^T \bar{w} - Y) = 0.$$

$$\bar{x}(\bar{x}^T \bar{w} - Y) + \lambda \bar{w} = 0$$

$$\bar{x}\bar{x}^T \bar{w} - \bar{x}Y + \lambda \bar{w} = 0.$$

$$(\bar{x}\bar{x}^T + \lambda I) \bar{w} = \bar{x}Y$$

$$\therefore \bar{w} = (\bar{x}\bar{x}^T + \lambda I)^{-1} \bar{x}Y \quad \left\{ \begin{array}{l} \text{pre multiply by} \\ (\bar{x}\bar{x}^T + \lambda I)^{-1} \end{array} \right\}$$

$$\boxed{\bar{w} = C^{-1}d}$$

(1.2) Now suppose we remove x_i from the training data, let c_i, d_i, \bar{w}_i be the corresponding matrices for removing x_i . Express c_i in terms of C and x_i . Express d_i in terms of d and x_i .

Ans. (i) We know: $C = \bar{x}\bar{x}^T + \lambda I$.

Removing $x_i \Rightarrow c_i = \bar{x}\bar{x}^T - \bar{x}_i \bar{x}_i^T + \lambda I$

$$\boxed{c_i = C - \bar{x}_i \bar{x}_i^T}$$

(ii) we know: $d = \bar{x}y$

$$\therefore d_i = (\bar{x} - \bar{x}_i) y_i = \bar{x}y - \bar{x}_i y_i = \boxed{d - \bar{x}_i y_i}$$

(1.3) Express c_i^{-1} in terms of c^{-1} and x_i .

Ans.

~~Using the~~

$$c_i = c - \bar{x}_i \bar{x}_i^T$$

$$\therefore c_i^{-1} = (c - \bar{x}_i \bar{x}_i^T)^{-1}$$

Using the Sherman-Morrisson formula:

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

$$\therefore c_i^{-1} = c^{-1} + \frac{c^{-1} \bar{x}_i \bar{x}_i^T c^{-1}}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}$$

(1.4) To show That: $\bar{w}_i = \bar{w} + (c^{-1} \bar{x}_i) \frac{-y_i + \bar{x}_i^T \bar{w}}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}$

Ans. We have: $\bar{w} = c^{-1}d$

$$\therefore \bar{w}_i = c_i^{-1}d_i$$

$$= \left(c^{-1} - \frac{c^{-1} \bar{x}_i \bar{x}_i^T c^{-1}}{1 - \bar{x}_i^T c^{-1} \bar{x}_i} \right) (\bar{x}_i \bar{y} - \bar{x}_i y_i)$$

$$= c^{-1}(\bar{x}_i \bar{y} - \bar{x}_i y_i) - \left(\frac{c^{-1} \bar{x}_i \bar{x}_i^T c^{-1}}{1 - \bar{x}_i^T c^{-1} \bar{x}_i} \right) (\bar{x}_i \bar{y} - \bar{x}_i y_i)$$

$$= c^{-1}d - c^{-1} \bar{x}_i y_i \cancel{c^{-1} \bar{x}_i \bar{y}} - \left(\frac{c^{-1} \bar{x}_i \bar{x}_i^T c^{-1}}{1 - \bar{x}_i^T c^{-1} \bar{x}_i} \right) (\bar{x}_i \bar{y} - \bar{x}_i y_i)$$

$$= -c^{-1} \bar{x}_i y_i + \bar{w} - \left(\frac{c^{-1} \bar{x}_i \bar{x}_i^T c^{-1}}{1 - \bar{x}_i^T c^{-1} \bar{x}_i} \right) (\bar{x}_i \bar{y} - \bar{x}_i y_i)$$

$$= \bar{w} - \frac{c^{-1} \bar{x}_i y_i (1 - \bar{x}_i^T c^{-1} \bar{x}_i) - c^{-1} \bar{x}_i \bar{x}_i^T c^{-1} (\bar{x}_i \bar{y} - \bar{x}_i y_i)}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}$$

$$= \bar{w} - \frac{c^{-1} \bar{x}_i y_i - \frac{1 - \bar{x}_i^T c^{-1} \bar{x}_i}{c^{-1} \bar{x}_i y_i \bar{x}_i^T c^{-1} \bar{x}_i} - c^{-1} \bar{x}_i \bar{x}_i^T c^{-1} (\bar{x}_i \bar{y} - \bar{x}_i y_i)}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}$$

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

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

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

$$\boxed{\bar{w}_i = \bar{w} + \frac{(c^{-1} \bar{x}_i) (-y_i + \bar{x}_i^T \bar{w})}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}}$$

(1.5) Show that: LOOCV error for removing the i^{th} training data is:

$$\bar{w}_i^T \bar{x}_i - y_i = \frac{\bar{w}^T \bar{x}_i - y_i}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}$$

Ans. $w_i = \bar{w} + \frac{(c^{-1} \bar{x}_i) (-y_i + \bar{x}_i^T \bar{w})}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}$

$$w_i^T = \bar{w}^T + \frac{(-y_i + \bar{w}^T \bar{x}_i) (\bar{x}_i^T (c^{-1})^T)}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}$$

$$w_i^T \bar{x}_i = \bar{w}^T \bar{x}_i + \frac{(-y_i + \bar{w}^T \bar{x}_i) (\bar{x}_i^T c^{-1} \bar{x}_i)}{1 - \bar{x}_i^T c^{-1} \bar{x}_i}$$

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

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

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

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

(1.6) Formula from 1.5:

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

- Time complexity for calculating: $C^{-1} = k^3$.
- Time taken to calc. $\bar{w}^T x_i - y_i = k$.
- Time taken to calc. $1 - \bar{x}_i^T C^{-1} \bar{x}_i = k^2$.

$$\therefore \text{Overall complexity} = \boxed{O(nk^2 + k^3)}$$

\rightarrow Usual way of computing Loocv: $\boxed{O(nk^3)}$

\therefore we can see that formula from (1.5) is faster.

Q.2 Naïve Bayes and Logistic Regression

(21) Give formula for computing $P(Y|X)$.

$Y \rightarrow \text{boolean}$

$X = (X_1, X_2)$, $X_1 \rightarrow \text{boolean}$; $X_2 \rightarrow \text{continuous}$

Ans. Let $P(X_1|Y=0) = \lambda_0^{X_1} (1-\lambda_0)^{1-X_1}$ {as X_1 is boolean}

$P(X_1|Y=1) = \lambda_1^{X_1} (1-\lambda_1)^{1-X_1}$

Also, let $P(X_2|Y=0) = \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(X_2-\mu_0)^2}{2\sigma_0^2}}$

$P(X_2|Y=1) = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(X_2-\mu_1)^2}{2\sigma_1^2}}$

We know that: $P(Y|X) = \frac{P(X|Y)P(Y)}{\sum_y P(X|Y=y)P(Y=y)}$

$$\therefore P(Y=1|X) = \frac{\lambda_1^{X_1} (1-\lambda_1)^{1-X_1} \cdot \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(X_2-\mu_1)^2}{2\sigma_1^2}} \cdot P(Y=1)}{\lambda_1^{X_1} (1-\lambda_1)^{1-X_1} \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(X_2-\mu_1)^2}{2\sigma_1^2}} \cdot P(Y=1) + \lambda_0^{X_1} (1-\lambda_0)^{1-X_1} \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(X_2-\mu_0)^2}{2\sigma_0^2}} \cdot P(Y=0)}$$

~~We can write for any (Y,y) as follows:~~

$$P(Y=0|X) = \frac{\lambda_0^{X_1} (1-\lambda_0)^{1-X_1} \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(X_2-\mu_0)^2}{2\sigma_0^2}} P(Y=0)}{\lambda_0^{X_1} (1-\lambda_0)^{1-X_1} \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(X_2-\mu_0)^2}{2\sigma_0^2}} P(Y=0) + \lambda_1^{X_1} (1-\lambda_1)^{1-X_1} \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(X_2-\mu_1)^2}{2\sigma_1^2}} P(Y=1)}$$

(2.2)

$$\begin{aligned}
 P(Y=1|X) &= \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)} \\
 &= \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}} \\
 &= \frac{1}{1 + \left[\exp \left(\log \left(\frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)} \right) \right) \right]} \\
 &= \frac{1}{1 + \left(\exp \left(\log \left(\frac{1-\theta}{\theta} \right) + \sum_i \log \left(\frac{P(X_i|Y=0)}{P(X_i|Y=1)} \right) \right) \right)}
 \end{aligned}$$

Consider following term:

$$\sum_i \log \left(\frac{P(X_i|Y=0)}{P(X_i|Y=1)} \right) = \sum_i \left[\log P(X_i|Y=0) - \log P(X_i|Y=1) \right]$$

let $\lambda_{i0} = P(X_i|Y=0)$, $\lambda_{i1} = P(X_i|Y=1)$ we can write $P(X_i|Y=0) = \lambda_{i0}^{x_i} (1-\lambda_{i0})^{1-x_i}$ and $P(X_i|Y=1) = \lambda_{i1}^{x_i} (1-\lambda_{i1})^{1-x_i}$

$$\begin{aligned}
 \therefore \sum_i [\log P(X_i|Y=0) - \log P(X_i|Y=1)] &= \sum_i \left[\log (\lambda_{i0}^{x_i} (1-\lambda_{i0})^{1-x_i}) - \log (\lambda_{i1}^{x_i} (1-\lambda_{i1})^{1-x_i}) \right] \\
 &= \sum_i \left\{ (x_i \log \lambda_{i0} + (1-x_i) \log (1-\lambda_{i0})) - (x_i \log \lambda_{i1} + (1-x_i) \log (1-\lambda_{i1})) \right\} \\
 &= \sum_i \left[x_i \log \lambda_{i0} + \log (1-\lambda_{i0}) - x_i \log (1-\lambda_{i0}) - x_i \log \lambda_{i1} - \log (1-\lambda_{i1}) + x_i \log (1-\lambda_{i1}) \right] \\
 &= \sum_i \left\{ x_i [\log \lambda_{i0} - \log (1-\lambda_{i0}) - \log \lambda_{i1} + \log (1-\lambda_{i1})] + \log (1-\lambda_{i0}) - \log (1-\lambda_{i1}) \right\}
 \end{aligned}$$

$$\therefore \sum_i^d \log \left(\frac{P(X_i|Y=0)}{P(X_i|Y=1)} \right) = \sum_i^d \theta_i X_i + \cancel{\sum_i^d \theta_i} \theta_{d+1}$$

$$\text{where : } \theta_i = \log \lambda_{i0} - \log(1-\lambda_{i0}) - \log \lambda_{i1} + \log(1-\lambda_{i1})$$

$$\theta_{d+1} = \sum_i^d \log(1-\lambda_{i0}) - \log(1-\lambda_{i1})$$

$$\therefore P(Y=1|X) = \frac{1}{1 + \exp\left(-\left(\sum_i^d \theta_i X_i + \theta_{d+1}\right)\right)}$$

Q.3 Implementation of SVM

QUESTION 3- IMPLEMENTATION OF SVMs

(3.1)

(1) Write svm dual objective as a quadratic program.
Write down $H, f, A, b, Aeq, beq, lb, ub$.

$$\begin{aligned} \max_{\alpha} \quad & \sum_{j=1}^n \alpha_j - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i \alpha_i y_j \alpha_j K(x_i, x_j) \\ \text{s.t.} \quad & \sum_{j=1}^n y_j \alpha_j = 0 \\ & 0 \leq \alpha_j \leq C \quad \forall j. \end{aligned}$$

Ans.

$$H = \text{diag}(Y) \cdot K(x_i, x_j) \cdot \text{diag}(Y)$$
$$\begin{bmatrix} y_{11} & 0 & \dots & 0 \\ 0 & y_{nn} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & y_{nn} \end{bmatrix}_{n \times n} \begin{bmatrix} x_1 x_1 & \dots & x_1 x_n \\ x_n x_1 & \dots & x_n x_n \end{bmatrix}_{n \times n} \begin{bmatrix} y_{11} & 0 & \dots & 0 \\ 0 & y_{nn} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & y_{nn} \end{bmatrix}_{n \times n}$$

$$f = \begin{bmatrix} -1 \\ -1 \\ \vdots \\ -1 \\ -1 \end{bmatrix}_{n \times 1}$$

$$A = [] \quad b = 0$$

$$Aeq = Y^T \quad [y_1 y_2 \dots y_n]_{1 \times n}$$

$$beq = 0$$

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

$$ub = C$$

$$\begin{bmatrix} C \\ \vdots \\ C \end{bmatrix}_{n \times 1}$$

(3.1.4) and (3.1.5)

For C = 0.1

a. Confusion Matrix =

168 16

3 180

b. Accuracy = 0.9482

c. Objective value of the SVM = 24.7648

d. Support Vectors = 339

For C = 10

a. Confusion Matrix =

178 6

4 179

b. Accuracy = 0.9728

c. Objective value = 112.1461

d. Support Vectors = 126

(3.2) Implement Multiclass SVM using Stochastic Gradient Descent

(3.2) Multiclass SVM using SGD—

(1) Subgradient of L_i wrt. w_{yi} :-

$$\frac{\partial L_i}{\partial w_{yi}} = \begin{cases} \frac{1}{n} w_{yi} - c x_i, \\ \frac{1}{n} w_{yi}, \end{cases}$$

if $w_{yi}^T x_i - w_{yi}^T x_i + 1 > 0$.
otherwise

(2) ~~Subgradient~~ Subgradient of L_i wrt $w_{\hat{y}_i}$:-

$$\frac{\partial L_i}{\partial w_{\hat{y}_i}} = \begin{cases} \frac{1}{n} w_{\hat{y}_i} + c x_i, \\ \frac{1}{n} w_{\hat{y}_i}, \end{cases}$$

if $w_{\hat{y}_i}^T x_i - w_{y_i}^T x_i + 1 > 0$
otherwise

(3) Subgradient of L_i wrt. w_j for $j \neq y_i$ & $j \neq \hat{y}_i$:-

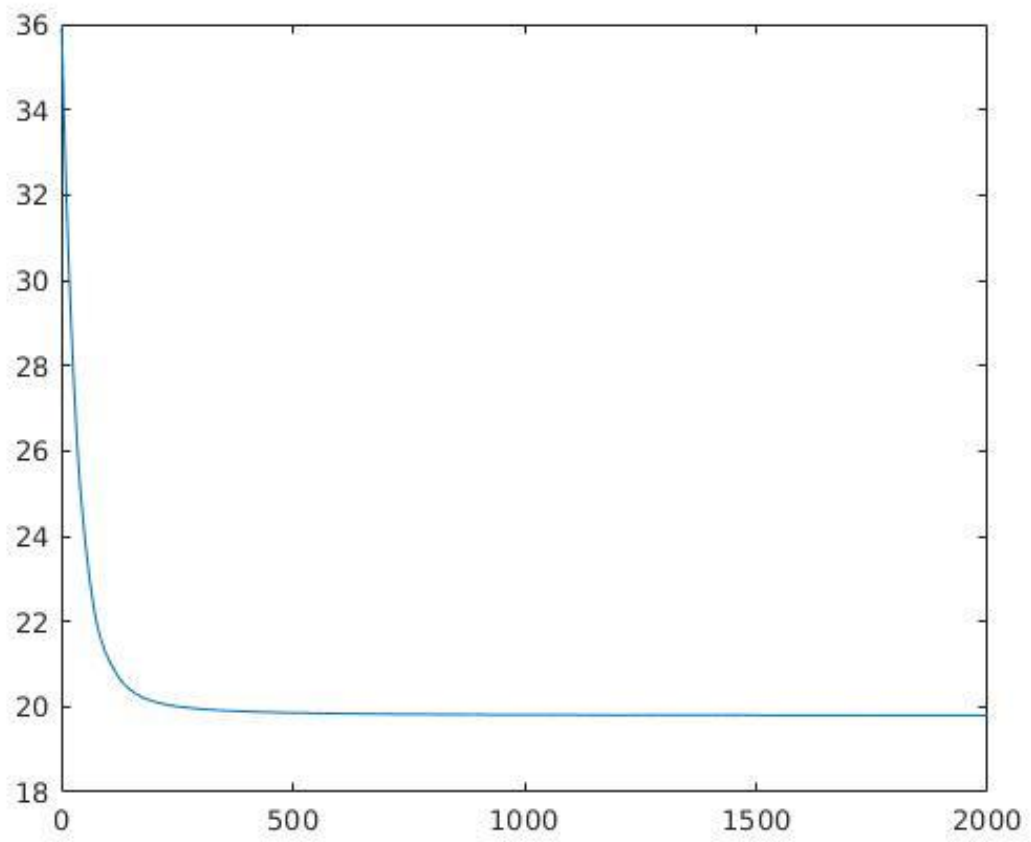
$$\frac{\partial L_i}{\partial w_j} = \frac{1}{n} w_j$$

where $j \neq y_i$ & $j \neq \hat{y}_i$

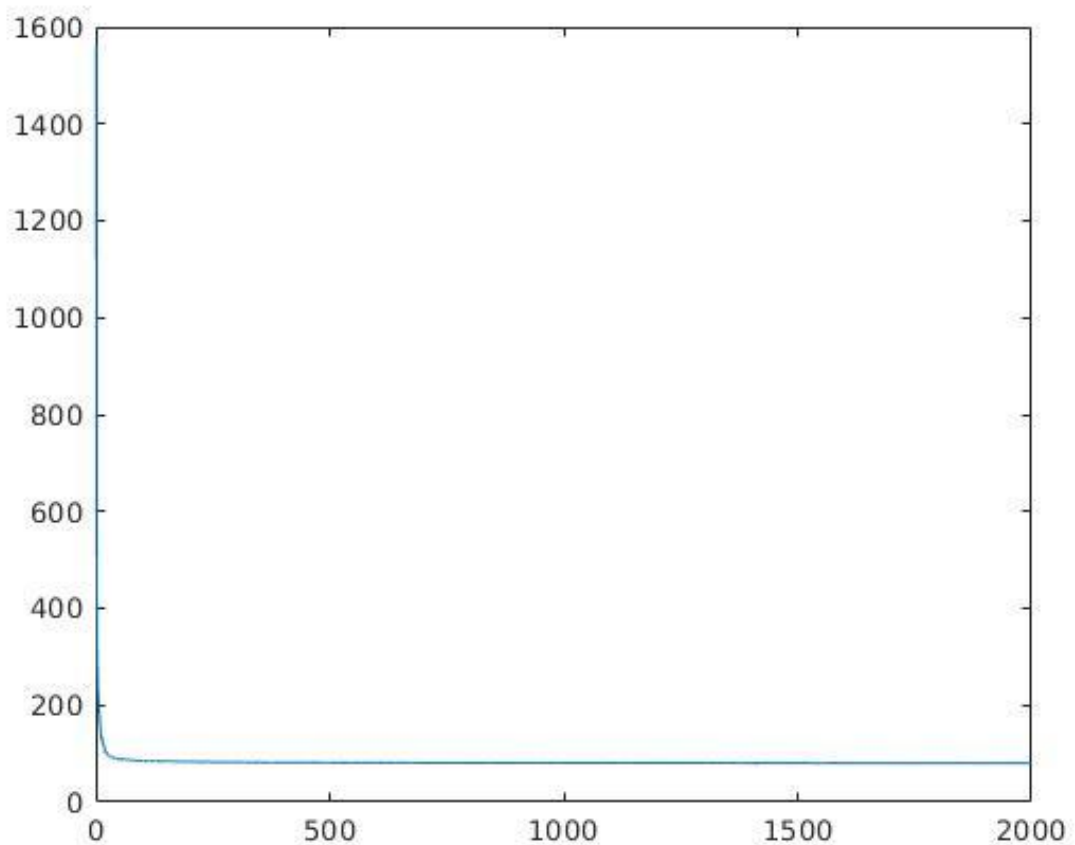
(3.2.5)

Using trD, trLb in q3 1 data.mat as your training set, run 2000 epochs over the dataset using $\eta_0 = 1$, $\eta_1 = 100$, $C = 0.1$ and $C = 10$. Plot the loss in Eq. (13) after each epoch. Compare with the objective value obtained in 3.1.4.

- For $C = 0.1$



- For $C = 10$



Objective Function Value:

For $C = 10$: 90

For $C = 0.1$: 22

(3.2.6) Using the W learned after 2000 epochs, report:

For $C = 0.1$

(a) Prediction error on $valD$, $valLb$, in `q3_1_data.mat` (test error) = 0.049

(b) The prediction error on trD , $trLb$ (training error) = 0.0304

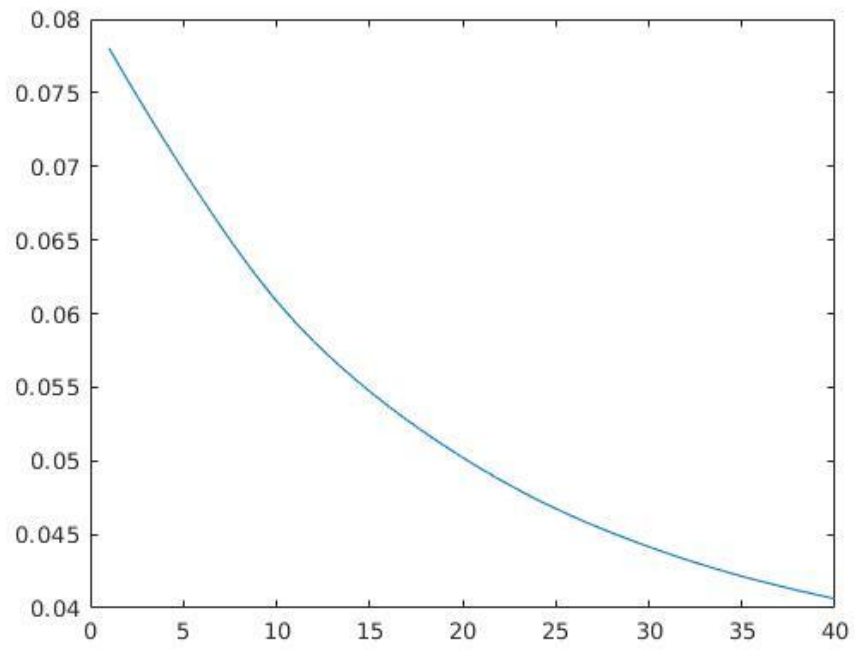
(c) 16.1101

For C = 10

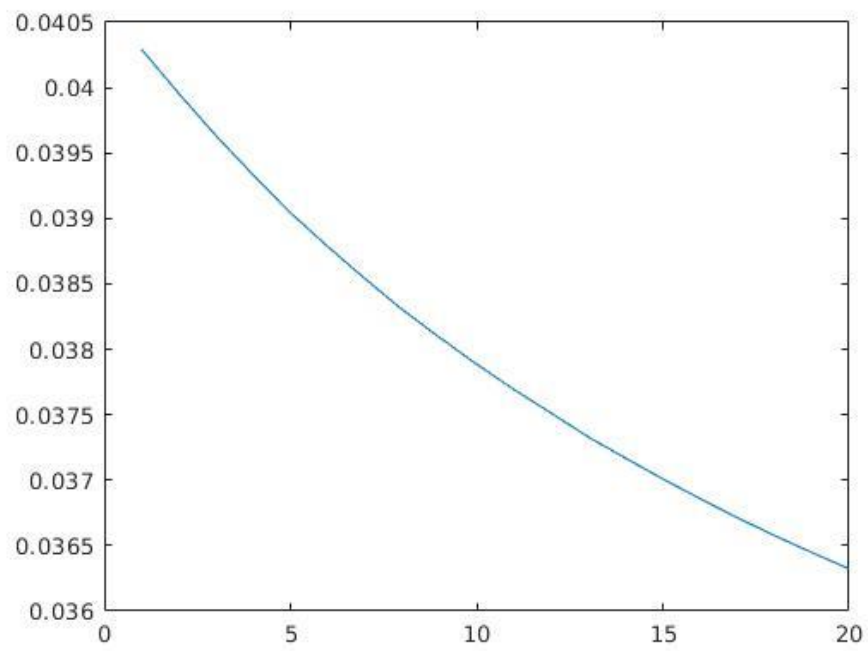
- (a) Prediction error on valD, valLb, in q3_1_data.mat (test error) = 0.0272
- (b) The prediction error on trD, trLb (training error) = 0
- (c) 120.95

(3.2.7) Kaggle

1. **Best accuracy:** 0.81097
2. **Parameters:**
 - a. C = 0.00001
 - b. Epochs = 40 (Training Data)
 - c. Epochs = 10 (Validation Data)
3. **Training Loss:**
 - a. **Training Data**



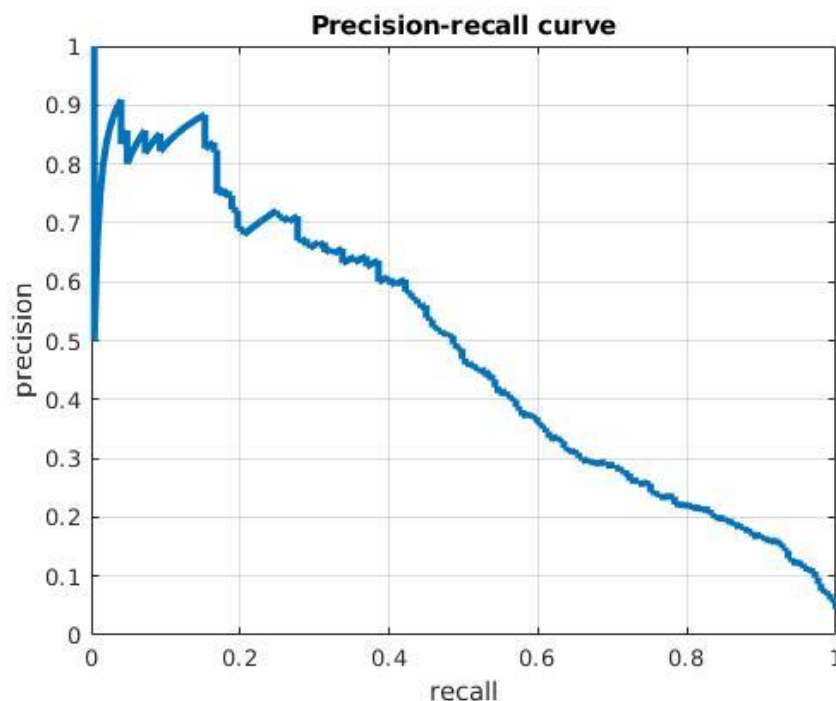
b. Validation Data



Q.4 SVM for object detection

(4.4.1) Use the training data in HW2 Utils.getPosAndRandomNeg() to train an SVM classifier. Use this classifier to generate a result file (use HW2 Utils.genRsltFile) for validation data. Use HW2 Utils.cmpAP to compute the AP and plot the precision recall curve. Submit your AP and precision recall curve (on validation data).

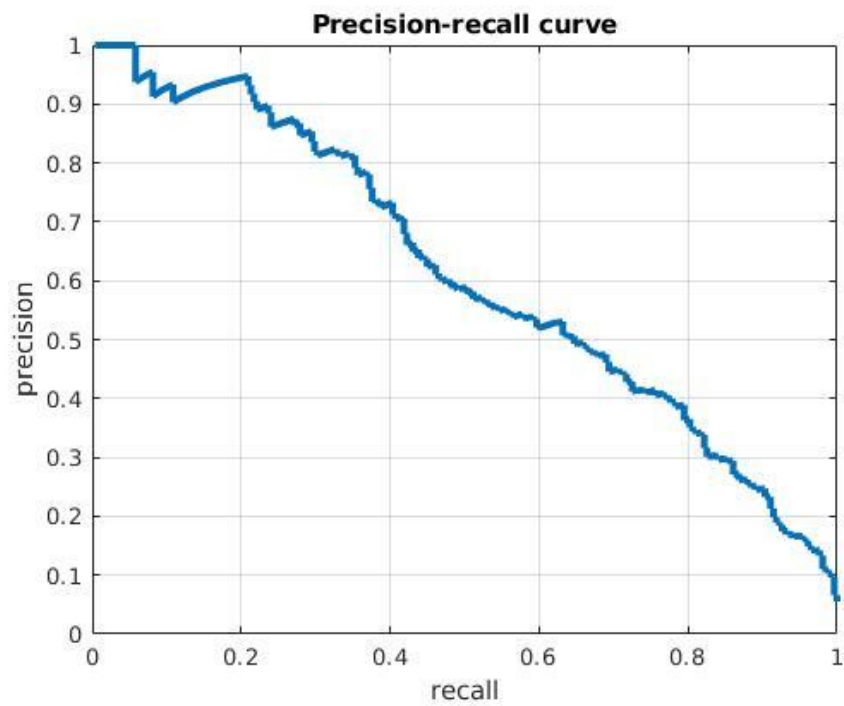
1. For $C = 10$
 - a. Precision Recall Curve



b. $ap = 0.4926$

2. For $C = 0.1$

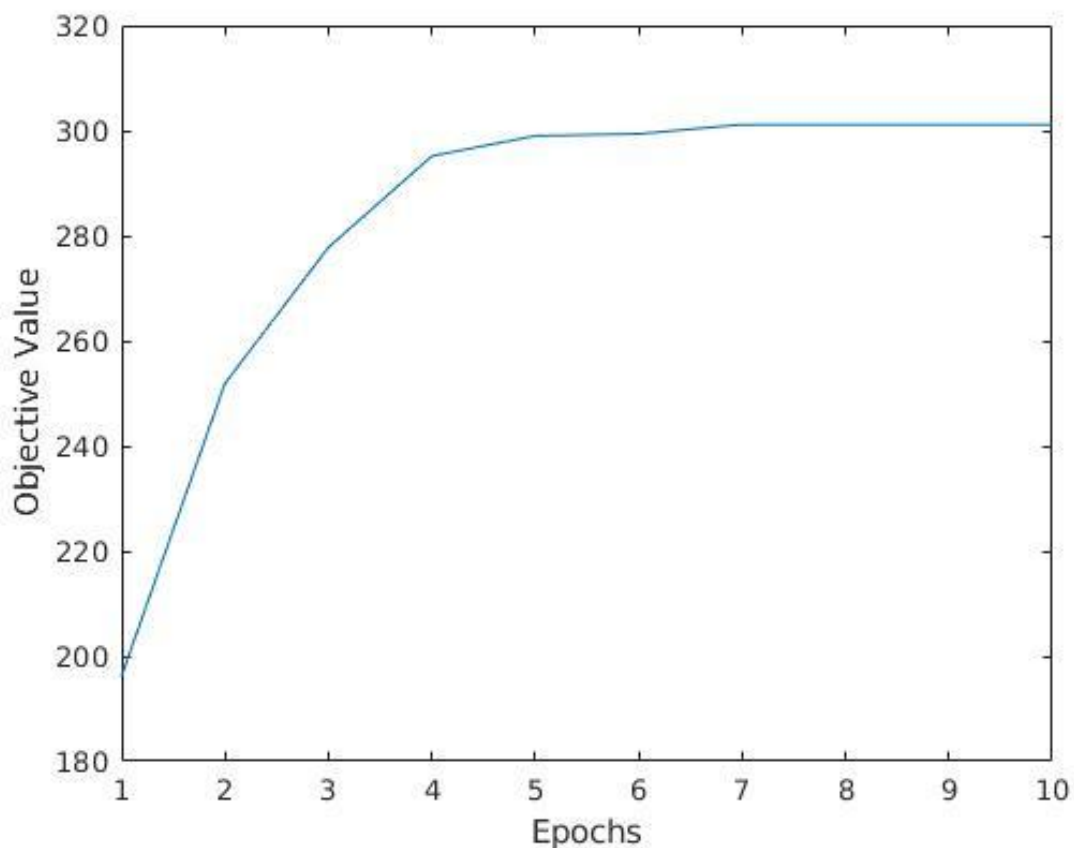
a. Precision Recall Curve



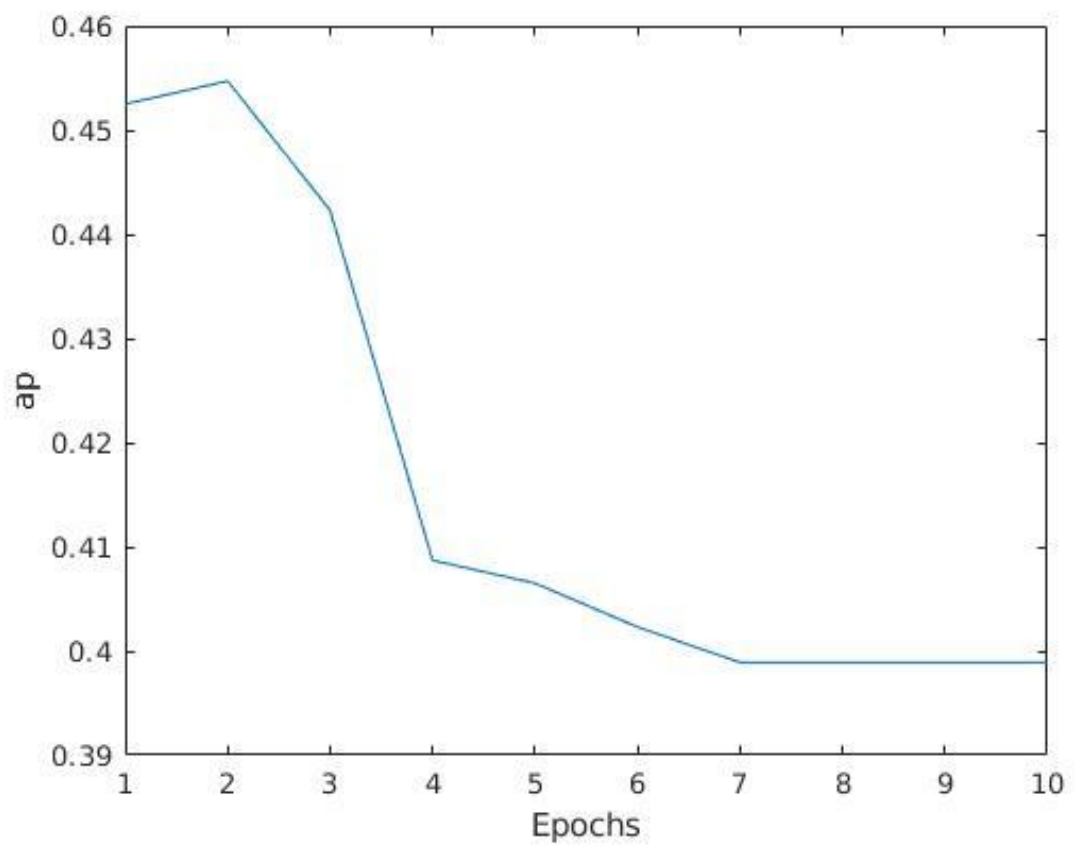
b. $ap = 0.6164$

(4.4.3) Run the negative mining for 10 iterations. Assume your computer is not so powerful and so you cannot add more than 1000 new negative training examples at each iteration. Record the objective values (on train data) and the APs (on validation data) through the iterations. Plot the objective values. Plot the APs.

1. Objective



2. AP



(4.4.4)

AP = 49.76