

CSE 512 - Machine Learning HW2

FAIZAAN CHARANIA - 111463646

Q1

The objective function for Ridge Regression is

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

1.1 Converting the above eqⁿ in vector form, we get

$$\min_{\bar{w}} \lambda \|\bar{w}\|^2 + \|\bar{X}^T \bar{w} - y\|^2 \quad - (ii)$$

To solve for \bar{w} , we take derivative of (ii)

w.r.t \bar{w} & equate to 0.

$$2\lambda \bar{w} + 2\bar{X}(\bar{X}^T \bar{w} - y) = 0$$

$$\lambda \bar{w} + \bar{X}(\bar{X}^T \bar{w} - y) = 0$$

$$\bar{X}y = (\bar{X}\bar{X}^T - \lambda I)\bar{w}$$

$$\bar{w} = (\bar{X}\bar{X}^T - \lambda I)^{-1} \bar{X}y$$

given $d = \bar{X}y$ & $c = (\bar{X}\bar{X}^T - \lambda I)$

~~$$\bar{w} = (\bar{X}\bar{X}^T - \lambda I)^{-1} \bar{X}y$$~~

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

Hence, proved.

1.2 Let $\mathcal{L}_i, d_i, \bar{w}_i$ correspond to matrices when we remove training data \bar{x}_i

For \mathcal{L}_i , each term will be missing a factor of \bar{x}_i which is $\bar{x}_i \bar{x}_i^T$

$$\therefore \mathcal{L}_i = C - \bar{x}_i \bar{x}_i^T \quad \text{---(iii)}$$

Similarly, for d_i we get

$$d_i = d - \bar{x}_i y_i \quad \text{---(iv)}$$

1.3 The Sherman-Morrison formula is

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

We know that

$$\mathcal{L}_i^{-1} = (\mathcal{L} - \bar{x}_i \bar{x}_i^T)^{-1} \quad (\text{from (iii)})$$

From (v), we get,

$$(\mathcal{L} - x_i x_i^T)^{-1} = \mathcal{L}^{-1} - \frac{\mathcal{L}^{-1} (-\bar{x}_i \bar{x}_i^T) \mathcal{L}^{-1}}{1 + \bar{x}_i^T \mathcal{L}^{-1} \bar{x}_i}$$

$$\therefore \mathcal{L}_i^{-1} = \mathcal{L}^{-1} + \frac{\mathcal{L}^{-1} \bar{x}_i \bar{x}_i^T \mathcal{L}^{-1}}{1 - \bar{x}_i^T \mathcal{L}^{-1} \bar{x}_i} \quad \text{--- (vi)}$$

Q1.4

We know that, $w = \mathcal{L}^{-1} d$

$$\therefore w_i = \mathcal{L}_i^{-1} d_i$$

\therefore from (vi) & (v)

$$\bar{w}_i = \mathcal{L}_i^{-1} d_i = \left[\mathcal{L}^{-1} + \frac{\mathcal{L}^{-1} \bar{x}_i \bar{x}_i^T \mathcal{L}^{-1}}{1 - \bar{x}_i^T \mathcal{L}^{-1} \bar{x}_i} \right] (d - x_i y_i)$$

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

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

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

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

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

$$\therefore \bar{w}_i = \bar{w} + C^{-1} \bar{x}_i \left[\frac{-y_i + \bar{x}_i^T \bar{w}}{1 - \bar{x}_i^T C^{-1} \bar{x}_i} \right] \quad \text{--- (11)}$$

Hence, proved

1.5 Leave one out error can be written as

$$\bar{x}_i^T \bar{\omega}_i - y_i \quad (\because \bar{\omega}^T \bar{x}_i = \bar{x}_i^T \bar{\omega})$$

from (ii) we know that

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

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

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

$$= \bar{\omega}^T \bar{x}_i - \cancel{\bar{\omega}^T \bar{x}_i \bar{x}_i^T C^{-1} \bar{x}_i} + \cancel{\bar{x}_i C^{-1} \bar{x}_i^T} (-y_i) + \cancel{\bar{x}_i C^{-1} \bar{x}_i^T} \bar{x}_i^T \bar{\omega} - y_i + \cancel{y_i \bar{x}_i^T C^{-1} \bar{x}_i}$$

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

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

$$\therefore \bar{x}_i^T \bar{\omega}_i - y_i = \frac{\bar{\omega}^T \bar{x}_i - y_i}{1 - \bar{x}_i^T C^{-1} \bar{x}_i}$$

1.6

In the previous question, we get

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

In this we need to calculate " C^{-1} " once, with complexity k^3 & we do the error calculation 'n' times with matrix multiplication complexity k^2

$$\therefore \text{Complexity} = nk^2 + k^3$$

In the usual way of computing we will have to per iteration complexity = k^3 (for inverse of C)

$$\therefore \text{Total complexity} = nk^3 \quad (\text{for } n \text{ iterations})$$

Q2

2.1 ~~N.B~~ N.B with Boolean & continuous variables

$$X \rightarrow Y$$

$$X = (X_1, X_2)$$

$$X_1 \sim \text{Bernoulli}(\theta_{x_i})$$

$$X_2 \sim \mathcal{N}(\mu_i, \sigma_i) \quad \forall i: i = \{0, 1\}$$

X_1 & X_2 have 3 parameters for each class (2 classes)

$$\therefore \text{num-params} = 3 \times 2 = 6 \quad (\text{Params for } X)$$

$$\begin{aligned} \text{Total Params} &= \text{Params for } X + \text{Param for } (Y) \\ &= 6 + 1 = 7 \end{aligned}$$

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

$$= \frac{P(X_1|Y) \cdot P(X_2|Y) \cdot P(Y)}{P(X_1|Y=0) \cdot P(X_2|Y=0) \cdot P(Y=0) + P(X_1|Y=1) \cdot P(X_2|Y=1) \cdot P(Y=1)}$$

$$P(X_1|Y=0) \cdot P(X_2|Y=0) \cdot P(Y=0) + P(X_1|Y=1) \cdot P(X_2|Y=1) \cdot P(Y=1)$$

For $Y=0$, we get,

$$P(Y=0|X) = \theta_{x_1 0} \times \frac{1}{\sigma_0 \sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \times (1 - \theta_y)$$

$$\left[\theta_{x_1 0} \times \frac{1}{\sigma_0 \sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \times (1 - \theta_y) \right] + \left[\theta_{x_1 1} \times \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \times \theta_y \right]$$

Q2

2.1 ~~N-B~~ N-B with Boolean & continuous variables

$$X \rightarrow Y$$

$$X = (X_1, X_2)$$

$$X_1 \sim \text{Bernoulli}(\theta_{x_i})$$

$$X_2 \sim \mathcal{N}(\mu_i, \sigma_i) \quad \forall_i i = \{0, 1\}$$

X_1 & X_2 have 3 parameters for each class (2 classes)

$$\therefore \text{num-params} = 3 \times 2 = 6 \quad (\text{Params for } X)$$

$$\begin{aligned} \text{Total Params} &= \text{Params for } X + \text{Param for } (Y) \\ &= 6 + 1 = 7 \end{aligned}$$

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

$$= \frac{P(X_1|Y) \cdot P(X_2|Y) \cdot P(Y)}{P(X_1|Y=0) \cdot P(X_2|Y=0) \cdot P(Y=0) + P(X_1|Y=1) \cdot P(X_2|Y=1) \cdot P(Y=1)}$$

$$P(X_1|Y=0) \cdot P(X_2|Y=0) \cdot P(Y=0) + P(X_1|Y=1) \cdot P(X_2|Y=1) \cdot P(Y=1)$$

For $Y=0$, we get,

$$P(Y=0|X) = \theta_{x_1 0} \times \frac{1}{\sigma_0 \sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \times (1-\theta_y)$$

$$\left[\theta_{x_1 0} \times \frac{1}{\sigma_0 \sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \times (1-\theta_y) \right] + \left[\theta_{x_1 1} \times \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \times \theta_y \right]$$

Similarly, we can find $P(Y=1|x)$

$$P(Y=1|x) = \theta_{x11} \times \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \times \theta_y$$

$$\left[\theta_{x11} \times \frac{1}{\sigma_1 \sqrt{2\pi}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \times \theta_y \right] + \left[\theta_{x10} \times \frac{1}{\sigma_0 \sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \times (1-\theta_y) \right]$$

2.2

$$X = (x_1, \dots, x_d)$$

$$X_i = \text{Bernoulli}(\theta_{ij}) \quad \forall i = \{1, \dots, d\} \quad \forall j = \{0, 1\}$$

$$Y = \text{Bernoulli}(\theta_y)$$

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

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

Consider,

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

$$= \theta_{ik}^{x_i} (1 - \theta_{ik})^{(1-x_i)}$$

$$P(x_i | Y=k) = \frac{1}{1 + \exp \left[\log \left(\frac{1 - \gamma_y}{\gamma_y} \right) + \sum \log \left(\frac{P(x_i | Y=0)}{P(x_i | Y=1)} \right) \right]}$$

$$= \frac{1}{1 + \exp \left[\log \left(\frac{1 - \gamma_y}{\gamma_y} \right) + \sum_{i=1}^d \log \left(\frac{\gamma_{i0}^{x_i} (1 - \gamma_{i0})^{(1-x_i)}}{\gamma_{i1}^{x_i} (1 - \gamma_{i1})^{(1-x_i)}} \right) \right]}$$

$$= \frac{1}{1 + \exp \left[\underbrace{\log \left(\frac{1 - \gamma_y}{\gamma_y} \right) + \sum_{i=1}^d \log \left(\frac{1 - \gamma_{i0}}{1 - \gamma_{i1}} \right)}_{-\theta_{d+1}} + \sum_{i=1}^d x_i \underbrace{\left[\log \frac{\gamma_{i0}}{\gamma_{i1}} - \log \left(\frac{1 - \gamma_{i0}}{1 - \gamma_{i1}} \right) \right]}_{-\theta_i} \right]}$$

$$\therefore P(Y=1 | x) = \frac{1}{1 + \exp \left(- \left(\sum_{i=1}^d \theta_i x_i + \theta_{d+1} \right) \right)}$$

$$\theta_{d+1} = - \left[\log \left(\frac{1 - \gamma_y}{\gamma_y} \right) + \sum_{i=1}^d \log \left(\frac{1 - \gamma_{i0}}{1 - \gamma_{i1}} \right) \right] = \log \left(\frac{\gamma_y}{1 - \gamma_y} \right) + \sum_{i=1}^d \log \left(\frac{1 - \gamma_{i1}}{1 - \gamma_{i0}} \right)$$

$$\theta_i = - \left[\log \left(\frac{\gamma_0}{\gamma_1} \right) - \log \left(\frac{1 - \gamma_{i0}}{1 - \gamma_{i1}} \right) \right] = \log \left(\frac{\gamma_1}{\gamma_0} \right) - \log \left(\frac{1 - \gamma_{i1}}{1 - \gamma_{i0}} \right)$$

Q3

3.1 Kernel SVM using quadprog

1. Based on the function definition of quadprog in MATLAB, we get

$$H = \text{diag}(1) \times \text{linear-kernel}(X) \cdot \text{diag}(1)$$

$$f = -1 \times \text{ones}(1, n)$$

$$A = []$$

$$b = []$$

} \because there are no inequality constraints

$$beq = 0$$

$$Aeq = Y^T$$

$$ub = (X \times \text{ones}(n, 1))$$

$$lb = \text{zeros}(n, 1)$$

where $n = \text{no. of samples}$
 $Y = \text{Training label vector } (n \times 1)$
 $X = \text{Training data matrix } (d \times n)$

$$\text{linear-kernel}(X) = X^T X$$

5) $\mathcal{L} = 10$

Validation Accuracy: 97.27

Objective value of SVM: 112.1461

No. of support vectors: 123

Confusion matrix = $\begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 178 & 6 \\ 4 & 179 \end{bmatrix}$

Q3.2

$$\mathcal{L}_i = \frac{1}{2n} \sum_{j=1}^k \|w_j\|^2 + C L(W, x_i, y_i)$$

1)

$$\frac{\partial \mathcal{L}_i}{\partial w_{y_i}} = \frac{1}{2n} \times 2 w_{y_i} + C \frac{\partial}{\partial w_{y_i}} L(W, x_i, y_i) \quad \text{--- (iii)}$$

$$\because L(W, x_i, y_i) = \max(w_{\hat{y}_i}^T x_i - w_{y_i} + 1, 0) \quad \text{--- (i x)}$$

where $\hat{y}_i = \underset{j \neq y_i}{\operatorname{argmax}} (w_j^T x_i)$
 Since (i x) is not differentiable,

we get

$$\frac{\partial}{\partial w_{y_i}} L(W, x_i, y_i)$$

$$\begin{aligned}
 \therefore \frac{\partial L_i}{\partial \omega_{\hat{y}_i}} &= \frac{1}{2n} \times 2 \omega_{\hat{y}_i} + C \times \frac{\partial}{\partial \omega_{\hat{y}_i}} L(W, x_i, y_i) \\
 &= \begin{cases} \omega_{\hat{y}_i} / n + C x_i & , \omega_{\hat{y}_i}^T x_i - \omega_{\hat{y}_i}^T x_{i+1} > 0 \\ \omega_{\hat{y}_i} / n & , \text{otherwise} \end{cases}
 \end{aligned}$$

3.

$$L_i = \frac{1}{2n} \sum_{j=1}^k \|\omega_j\|^2 + C \times L(W, x_i, y_i)$$

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

$$\frac{\partial L_i}{\partial \omega_j} = \frac{1}{2n} \times 2 \omega_j + C \times \frac{\partial}{\partial \omega_j} L(W, x_i, y_i)$$

$\therefore L(W, x_i, y_i)$ is independent of ω_j (from (ix))

$$\frac{\partial L(W, x_i, y_i)}{\partial \omega_j} = 0$$

$$\therefore \frac{\partial L_i}{\partial \omega_j} = \frac{\omega_j}{n}$$

4. code submitted

5. Objective value
3.1

3.2

$C = 0.1$

24.76

19.802

$C = 10$

112.14

80.62

Plots attached. (Compared to the value is 3.1.4
SGD performs marginally better in t)^x

6.

6.

Considering $C = 10$ (Assuming prediction error is loss ^{or} accuracy)

	Assume	Hinge Loss	Accuracy
a)	Validation	56.14	97.27
b)	Train	2.01	100

$$c) \sum_{j=1}^k \|w_j\|^2 = \cancel{50.64} 16.10$$

$$7) \text{ Validation accuracy} = 80.33\%$$

$$\text{Test accuracy} = 73.2\%$$

Parameters

$$\text{epochs} = 30$$

$$\cancel{\text{eta}} \eta_0 = 0.05$$

$$\eta_1 = 100$$

$$C = 0.07$$

$$\text{batch size} = 20$$

(Plot Attached)

Q4

~~4.4.1~~ AP = 0.5785 (validation data)

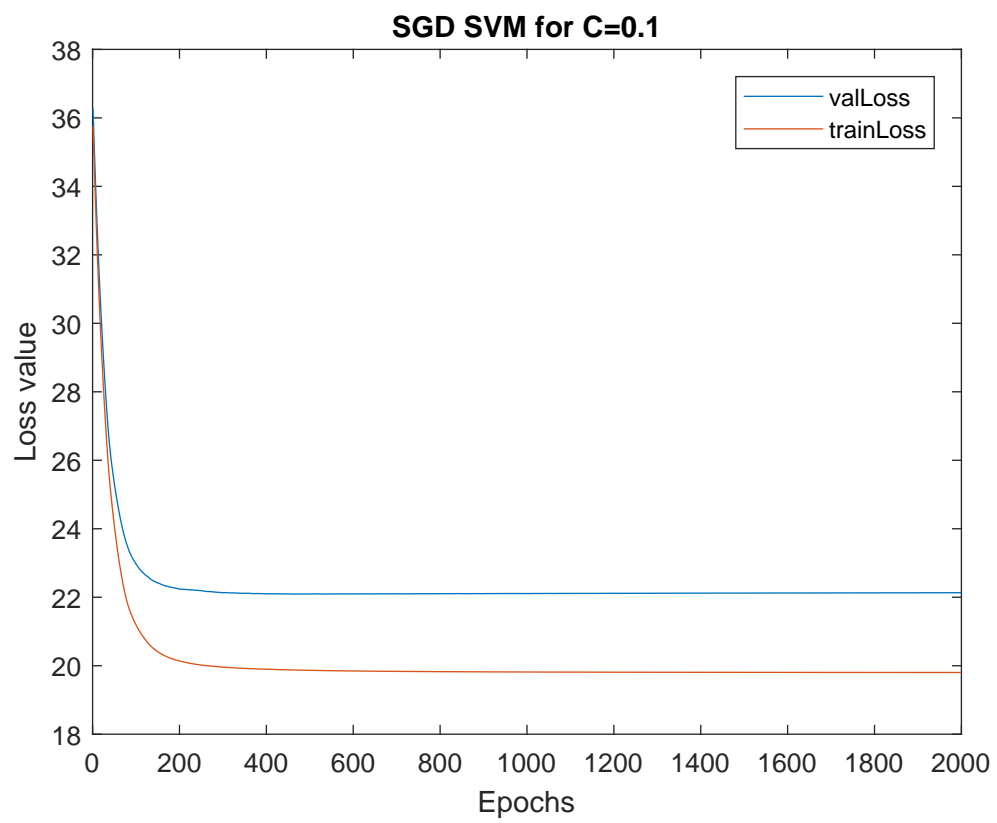
Precision-recall plot attached.

4.4.2 code submitted

4.4.3 Plots attached

4.4.4 AP% = 84.90%

SBU ID \Rightarrow 111463646



SGD SVM for C=10

