# EECS 4404

## Assignment 2

**Name: Bochao Wang**

**Student ID: 215237902**

**Prism: bochao**

**Date: March. 14th**

## 1. Gradient computation

Let $X \in \mathbb{R}^{d*d}$ be some matrix. Consider the function $f : \mathbb{R}^d \to \mathbb{R}$ defined by $f(w) = w^T X w$, show that the gradient of this function with repect to $w$ is $\triangledown f(w) = w^T(X + X^T)$

- Proof:

  ○ $\triangledown f(w) = \left( \frac{\partial f}{\partial w_1}, \quad \frac{\partial f}{\partial w_2}, \quad \cdots \quad \frac{\partial f}{\partial w_d} \right)$

  ○ $f(w) = \begin{bmatrix} w_1 & w_2 & w_3 & \cdots & w_d \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \cdots & x_{1d} \\ x_{21} & x_{22} & x_{23} & \cdots & x_{2d} \\ x_{31} & x_{32} & x_{33} & \cdots & x_{3d} \\ \vdots & & & & \vdots \\ x_{d1} & x_{d2} & x_{d3} & \cdots & x_{dd} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{bmatrix}$

  ○ $f(w) = \begin{bmatrix} \sum_{i=1}^{d} w_i x_{i1}, & \sum_{i=1}^{d} w_i x_{i2}, & \cdots & \sum_{i=1}^{d} w_i x_{id} \end{bmatrix} w = \sum_{j=1}^{d} \sum_{i=1}^{d} w_i x_{ij} w_j$

  ○ Simplify the function $f(w) = \sum_{j=1}^{d} (x_{jj} w_j^2 + \sum_{i \neq j} w_i x_{ij} w_j)$

  ○ we can compute the partial derivative with element $k$,

  $\frac{\partial}{\partial w_k} [\sum_{j=1}^{d} (x_{jj} w_j^2 + \sum_{i \neq j} w_i x_{ij} w_j)] = \sum_{j=1}^{d} (\frac{\partial}{\partial w_k} x_{jj} w_j^2 + \sum_{i \neq j} \frac{\partial}{\partial w_k} w_i x_{ij} w_j) = 2 x_{kk} w_k + \sum_{i \neq k} x_{ki} w_i + \sum_{i \neq k} w_i x_{ik}$

  ○ $\because x_{kk} w_{kk} = x_{kk} w_{kk}$ Then, $\frac{\partial}{\partial w_k} [\sum_{j=1}^{d} (x_{jj} w_j^2 + \sum_{i \neq j} w_i x_{ij} w_j)] = \sum_{i=1}^{d} x_{ki} w_i + \sum_{i=1}^{d} w_i x_{ik}$

  ○ $\triangledown f(w) = \begin{bmatrix} \sum_{i=1}^{d} x_{1i} w_i + \sum_{i=1}^{d} w_i x_{i1}, & \sum_{i=1}^{d} x_{2i} w_i + \sum_{i=1}^{d} w_i x_{i2}, & \cdots & \sum_{i=1}^{d} x_{di} w_i + \sum_{i=1}^{d} w_i x_{id} \end{bmatrix}$

  ○ $\triangledown f(w) = \begin{bmatrix} \sum_{i=1}^{d} w_i x_{i1}, & \sum_{i=1}^{d} w_i x_{i2}, & \cdots & \sum_{i=1}^{d} w_i x_{id} \end{bmatrix} + \begin{bmatrix} \sum_{i=1}^{d} x_{1i} w_i, & \sum_{i=1}^{d} x_{2i} w_i, & \cdots & \sum_{i=1}^{d} x_{di} w_i \end{bmatrix}$

  ○ $\begin{bmatrix} \sum_{i=1}^{d} w_i x_{i1}, & \sum_{i=1}^{d} w_i x_{i2}, & \cdots & \sum_{i=1}^{d} w_i x_{id} \end{bmatrix} = w^T X$

  ○ $\begin{bmatrix} \sum_{i=1}^{d} x_{1i} w_i, & \sum_{i=1}^{d} x_{2i} w_i, & \cdots & \sum_{i=1}^{d} x_{di} w_i \end{bmatrix} = w^T X^T$

  ○ Thus, $\triangledown f(w) = w^T X + w^T X^T = w^T(X + X^T)$

## 2. Stochastic Gradient Descent

Recall the logistic loos function, point wise defined as

$l^{logist}(y_{\boldsymbol{w}}, (\boldsymbol{x}, t)) = ln(1 + exp(-t < \boldsymbol{w}, \boldsymbol{x} >))$

We know that empirical logistic loss over a dataset

$\mathcal{L}^{logist}(y_{\boldsymbol{w}}) = \frac{1}{N} \sum_{n=1}^{N} ln(1 + exp(-t_n < \boldsymbol{w}, \boldsymbol{x}_n >))$

(a) Compute the gradient with respect to $w$ of the logistic loss $l^{logist}(y_{\boldsymbol{w}}, (\boldsymbol{x}, t))$ on a data point $(x, t)$.

- Solve:

  ○ $l^{logist}(y_{\boldsymbol{w}}, (\boldsymbol{x}, t)) = ln(1 + exp(-t * \sum_i^N w_i x_i))$

  ○ $\triangledown f(\boldsymbol{w}) = \left( \frac{\partial f}{\partial w_1}, \quad \frac{\partial f}{\partial w_2}, \quad \cdots \quad \frac{\partial f}{\partial w_d} \right)$

  ○ $\triangledown l^{logist}(y_{\boldsymbol{w}}, (\boldsymbol{x}, t)) = \left( \frac{1}{1 + exp(-t * \sum_i^N w_i x_i)} \frac{\partial(1 + exp(-t * \sum_i^N w_i x_i))}{\partial w_1}, \quad \frac{1}{1 + exp(-t * \sum_i^N w_i x_i)} \frac{\partial(1 + exp(-t * \sum_i^N w_i x_i))}{\partial w_2}, \quad \cdots, \quad \frac{1}{1 + exp(-t * \sum_i^N w_i x_i)} \frac{\partial(1 + exp(-t \cdots}{\partial u} \right.$

  ○ we compute the $k$ element:

  ○ $\frac{\partial(1 + exp(-t \sum_i^N w_i x_i))}{\partial w_k} = e^{-t \sum_i^N w_i x_i} \frac{\partial(-t \sum_i^N w_i x_i)}{\partial w_k} = e^{-t \sum_i^N w_i x_i}(-t x_k) = e^{-t < \boldsymbol{w}, \boldsymbol{x} >}(-t x_k)$

  ○ $\triangledown l^{logist}(y_{\boldsymbol{w}}, (\boldsymbol{x}, t)) = \left( \frac{e^{-t < \boldsymbol{w}, \boldsymbol{x} >}}{1 + exp(-t * \sum_i^N w_i x_i)}(-t x_1), \quad \frac{e^{-t < \boldsymbol{w}, \boldsymbol{x} >}}{1 + exp(-t * \sum_i^N w_i x_i)}(-t x_2), \quad \cdots, \quad \frac{e^{-t < \boldsymbol{w}, \boldsymbol{x} >}}{1 + exp(-t * \sum_i^N w_i x_i)}(-t x_d) \right)$

  ○ $\triangledown l^{logist}(y_{\boldsymbol{w}}, (\boldsymbol{x}, t)) = \left( \frac{e^{-t < \boldsymbol{w}, \boldsymbol{x} >}}{1 + exp(-t < \boldsymbol{w}, \boldsymbol{x} >)}(-t x_1), \quad \frac{e^{-t < \boldsymbol{w}, \boldsymbol{x} >}}{1 + exp(-t < \boldsymbol{w}, \boldsymbol{x} >)}(-t x_2), \quad \cdots, \quad \frac{e^{-t < \boldsymbol{w}, \boldsymbol{x} >}}{1 + exp(-t < \boldsymbol{w}, \boldsymbol{x} >)}(-t x_d) \right)$

  ○ $\triangledown l^{logist}(y_{\boldsymbol{w}}, (\boldsymbol{x}, t)) = \frac{e^{-t < \boldsymbol{w}, \boldsymbol{x} >}}{1 + exp(-t < \boldsymbol{w}, \boldsymbol{x} >)}(-t \boldsymbol{x})$

(b) Describe Stochastic Gradient Descent with a fixed stepsize η with respect to the logistic loss

- Solve:

  1. $In : data\ D = ((\boldsymbol{x}_1, t_1), (\boldsymbol{x}_2, t_2), \ldots, (\boldsymbol{x}_N, t_N))$
  2. $Parameters : \eta,\ T,\ n$
  3. $Initialize : \boldsymbol{w}_0 = 0$
  4. $For\ i = 0, \ldots, T$
  5. Choose a random data point $(\boldsymbol{x}^*, t^*)$ form $D$
  6. so that $\boldsymbol{w}_{i+1} = \boldsymbol{w}_i - \eta * \frac{e^{-t^* < \boldsymbol{w}_i, \boldsymbol{x}^* >}}{1 + exp(-t^* < \boldsymbol{w}_i, \boldsymbol{x}^* >)}(-t^* \boldsymbol{x}^*)$
  7. $Output : \boldsymbol{w} = \frac{1}{n} \sum_{k=T-n}^{T} \boldsymbol{w}^k$

  ○ $T$ is the number of updates
  ○ $\eta$ is fixed stepsize

- $n$ compute average of last $n$ iteration

(c) Compare the derived algorithm to the SGD algorithm we derived for SVM. What is similar? What is different? On what type of data would you expect to see different behavior?

- Solve:
  - Similar :
    - The part (b) algrothim and we derived for SVM all use the gradient descent method to update and find the optimal $w$ so that minimize the loss.
    - Both algorithms would sample a data point $(w_i, t_i)$ uniformly at the random from dataset D to compute the gradient to update
    - Both of them are convex optimal problems
    - Both algorithms compute the average $w$ of last $n$ iterations.
    - Both of them would be initialized $w_0 = 0$
  - Different:
    - part (b) use fixed stepsize $\eta$, SVM we derived use stepsize sequence $\eta * j = \frac{1}{\lambda * j}$ where $j$ is current step number. It means in SVM, the stepsize would less and less in each step, it would more accuracy to be closed minimum.
    - Part(b) consider the gradient of logistic loss, because logistic loss would be differenviable in every point, the gradient and subgtadient would be same. However, SVM consider the subgradient of hinge loss, because the point at $t_i < w, x_i >$ would not be differentiable. The logistic loss would close to 0, but never be 0; the hinge loss would be 0 if $t_i < w, x_i > \ > 1$, it means that the suport vectors would determine the $w$, but Part(b) algrothim every points, which classify correctlly, also update the $w$.
    - soft SVM allows sparse support vectors. It means soft SVM would tolerant the missclassified point in the margin.
    - The gradient of hinge loss function is a piecewise-defined function, but the gradient of logistic loss is comtinuous.
  - The subgradient of SVM is 0, if $t_i < w, x_i > \ > 1$, it means if the point is far from the margin and classified correctly, these points woud not update the $w$, but the SGD of logistic loss would always update the $w$. Thus, if the dataset contains most points where are in the $t_i < w, x_i > \ > 1$ i.e. almost point would far from the margin and classify correctly. In such dataset, SVM would update very few steps, but SGD of logistics loss would always update.

### 3. SoftSVM optimization

(a) Implement SGD for SoftSVM.

`soft_svm.m`

```matlab
function [w, hi_loss, bi_loss] = soft_svm(D,T,lambda, n)
% if n=0, only print last step w, otherwise print average of
% last n iteration

% Initialize
[N, clo] = size(D);
x = D(:,1:clo-1); % features
t = D(:,clo); % label {-,+}
theta=zeros(1,clo-1);

hi_loss = zeros(T,1);
bi_loss = zeros(T,1);

% store `w` in each iteration
ws=zeros(T,clo-1);

w_j = (1/(lambda))*theta; % w_0

for j = 1:T
    % Choose the i in uniform distribution
    i = unidrnd(N);

    if t(i) * dot(w_j, x(i,:)) < 1

        theta = theta + t(i)*(x(i,:));

    end

    % update `w`
    w_j = (1/(j*lambda))*theta;
```

```
32        % store the `w` in `ws`
33          ws(j,:) = w_j;
34
35        % Track the emprical and hinge loss
36          hi_loss(j,:) = emp_loss(w_j, D, 'hinge');
37          bi_loss(j,:) = emp_loss(w_j, D, 'binary');
38    end
39
40    if(n ~= 0)
41        w = (1/n) * sum(ws((T-n:T),:));
42    else
43        w = w_j;
44    end
```

emp_loss.m

```
1    function emr_loss = emp_loss(w, D, flag)
2    % L_D(w)=(1/N)sum(loss(w,(x_i, t_i)))
3    [N,clo] = size(D);
4    l=0;
5    for i = 1:N
6        x_i = D(i,1:clo-1);
7        t_i = D(i,clo);
8        if(strcmpi(flag, 'hinge'))
9            l = l + hinge_loss(w, x_i, t_i);
10       elseif(strcmpi(flag, 'binary'))
11           l = l + binary_loss(w, x_i, t_i);
12       end
13   end
14   emr_loss = (1/N) * l;
```

hinge_loss.m

```
1    function h_loss = hinge_loss(w, x_i, t_i)
2    % hinge_loss(w,(x,t)) = max{0, 1-t<w,x>}
3    if t_i * dot(w,x_i) > 1
4        h_loss = 0;
5    else
6        h_loss = 1 - t_i * dot(w,x_i);
7    end
```

binary_loss.m

```
1    function b_loss = binary_loss(w, x_i, t_i)
2    % binary_loss(w,(x,t)) = 0[y(x) = t]
3    if t_i * dot(w,x_i) >= 0
4        b_loss = 0;
5    else
6        b_loss = 1;
7    end
```

(b) Run the optimization method with various values for the regularization parameter $\lambda = \{100, 10, 1, .1, .01, .001\}$ on the data.
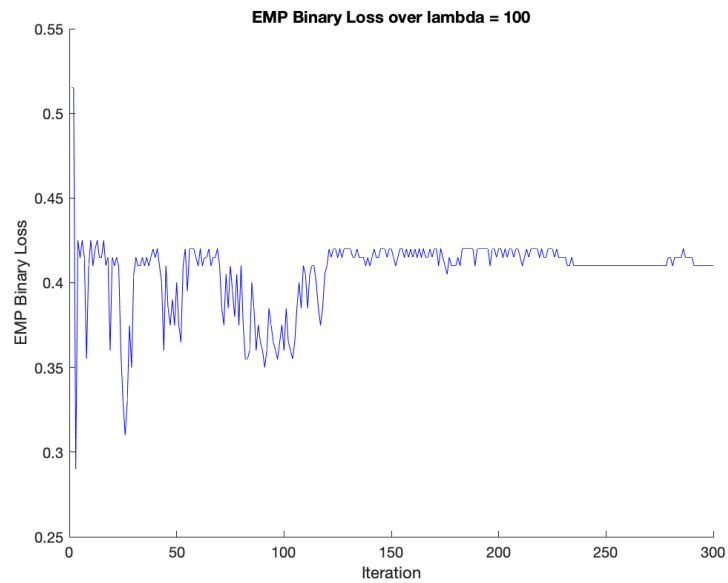
main.m

```
1    % PART (b) Run with lambda = {100, 10, 1, .1, .01, .001} on the data
2    % and print empiral hinge loss and empirical binary loss
3    clc;clear;close all;
4    dataset = load_data("bg.txt");
5
6    % add bias
7    [N,C] = size(dataset);
8    bias = ones(N,1);
```
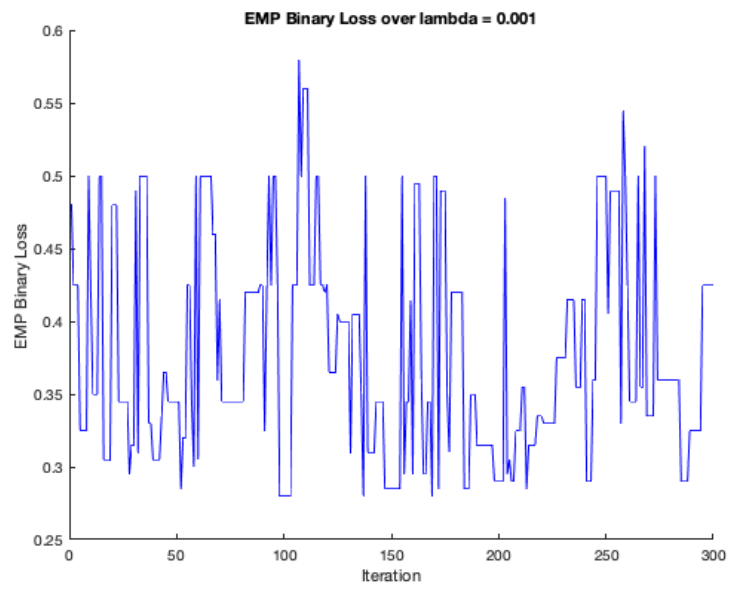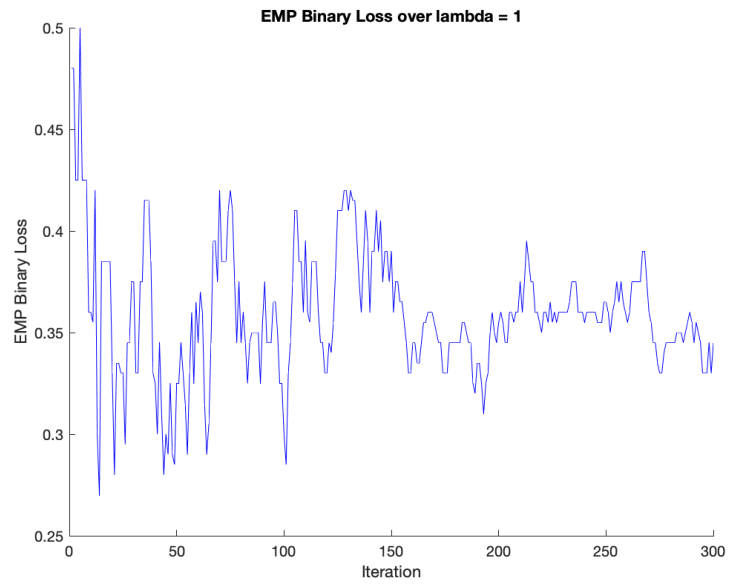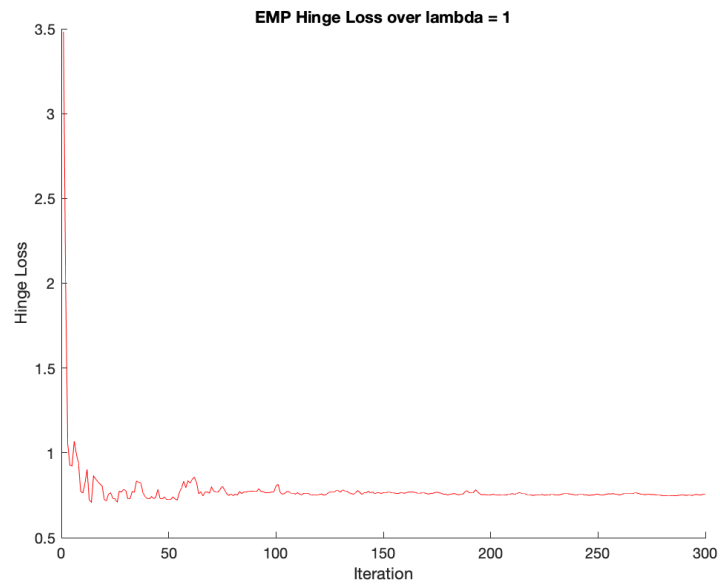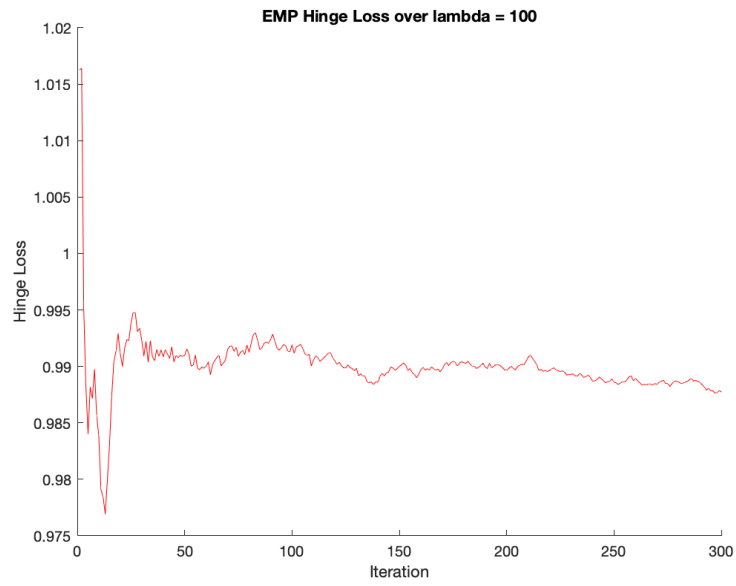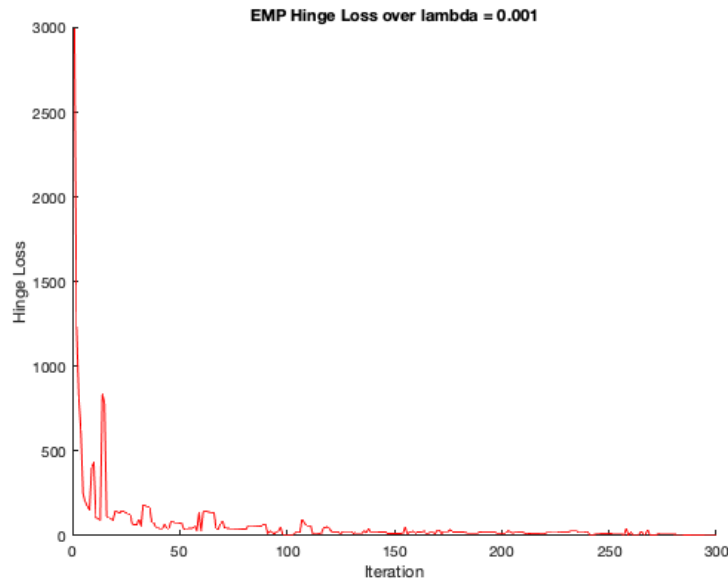
```matlab
 9  dataset = [dataset(:,1:C-1) bias dataset(:,C)];
10
11  num_updates = 300;
12  s_lambda = [100, 10, 1, .1, .01, .001];
13  n = 1;
14
15  for i = 1:length(s_lambda)
16      [w,hi_loss,bi_loss] = soft_svm(dataset, num_updates, s_lambda(i), n);
17
18      figure
19      title(['EMP Hinge Loss over lambda = ', num2str(s_lambda(i))])
20      hold on
21  %     set(gca, 'YScale', 'log'); % set y-axis as log
22      plot(hi_loss, 'Color', 'red');
23      xlabel('Iteration');
24      ylabel('Hinge Loss');
25      hold off
26
27      figure
28      title(['EMP Binary Loss over lambda = ', num2str(s_lambda(i))])
29      hold on
30  %     set(gca, 'YScale', 'log'); % set y-axis as log
31      plot(bi_loss, 'Color', 'blue');
32      xlabel('Iteration');
33      ylabel('EMP Binary Loss');
34      hold off
35  end
36
37  fprintf("part(b): Finished! Please press enter to continue!\n")
38  pause;
39  clc;close all;
```



EMP Binary Loss over lambda = 100

**EMP Binary Loss over lambda = 1**

**EMP Binary Loss over lambda = 0.001**

**EMP Hinge Loss over lambda = 100**

**EMP Hinge Loss over lambda = 1**

**EMP Hinge Loss over lambda = 0.001**

(c) Discuss the plots. Are the curves monotone? Are they approximately monotone? Why or why not? How does the choice of $\lambda$ affect the optimization? How would you go about finding a linear predictor of minimal binary loss?

- Solve:
  - No, the curves are not strictly monotone. Because in each iteration, the SGD for SoftSVM pick random data point to compute the subgradient, so the loss may not decreasing. For example, sometimes, it may pick a noice data point, the loss would larger than previous iteration; sometimes, it may pick a same point as previous iteration, the loss would very small. Thus, it not monotone.
  - Yes, they are approximately monotone. The trend of loss would montonely decreasing. Because gradient method would always give the direction which decreasing fast. The loss would decrease globally.
  - Because stepsize sequence $\eta^j = \frac{1}{\lambda * j}$ and the $\lambda$ is on the denominator, when $\lambda$ is larger, the stepsize would be small, and multiply the iteration number, the stepsize would smaller and smaller, it would close to the minimum slower, and it may need more iterations to close the minimum points . However, if $\lambda$ is very small, stepsize would very big initially, so it would overstep, it may just go across the optimal solution and increase the loss. Also, if $\lambda$ is very larger, the function would be simpler.
  - Because minimize the binary loss computation would very hard, we use its upper bounded loss i.e. hinge loss to substitute. To find a linear predictor of minimal binary loss, we use soft SVM and track the binary loss. we need try and modify the several different $\lambda$ and $T$ the minimum loss and obtain the weight in the SVM would also for linear predictor of minimal binary loss.

(e-f) Multi-class predictor

`main.m`

```matlab
1   % PART (d-e) Spilt the data set and train three binary linear predictors
2   clc;clear;close all;
3
4   % load the dataset and spilt
5   D = load('seeds_dataset.txt');
6   % + bias
7   [N,C] = size(D);
8   bias = ones(N,1);
9   D = [D(:,1:C-1) bias D(:,C)];
10
11  [N,C] = size(D);
12  [D_1, D_2, D_3] = spiltDateset(D); % split and add bias cloumn
13  num_updates = 10000;
14  lambda = 0.01;
15  n = 1;
16  binary_losses=zeros(3,1);
17  % ws = zeros(3,C-1);
18  ws = zeros(3,C-1); % + bias
19
```

```
20
21   min_loss = N;
22   min_ws_loss = zeros(3,1);
23
24
25   for i = 1:50
26
27       [w_1,~,bi_loss_1] = soft_svm(D_1, num_updates, lambda, n);
28       binary_losses(1) = emp_loss(w_1, D_1, 'binary');
29       ws(1,:) = w_1;
30
31       [w_2,~,bi_loss_2] = soft_svm(D_2, num_updates, lambda, n);
32       binary_losses(2) = emp_loss(w_2, D_2, 'binary');
33       ws(2,:) = w_2;
34
35       [w_3,~,bi_loss_3] = soft_svm(D_3, num_updates, lambda, n);
36       binary_losses(3) = emp_loss(w_3, D_3, 'binary');
37       ws(3,:) = w_3;
38
39       DR = D(:,C);
40       for j = 1:N
41           x_j = D(j, 1:C-1);
42           t_j = D(j,C);
43           [~, I] = max([dot(ws(1,:),x_j), dot(ws(2,:),x_j), dot(ws(3,:),x_j)]);
44           DR(j,1) = I;
45       end
46       Diff = DR(:,1) ~= D(:,C);
47       loss = sum(Diff)
48
49       if (loss < min_loss)
50           min_loss = loss;
51           min_ws_loss = binary_losses;
52           opt_DR = DR;
53           opt_ws = ws;
54       end
55   end
```

`spiltDateset.m`

```
1    function [D_1, D_2, D_3] = spiltDateset(D)
2    [~,C] = size(D);
3    % each spilt D into two part {-1,1}
4    D_1 = D;
5    D_1(D(:,C)==1, C) = 1;
6    D_1(D(:,C)~=1, C) = -1;
7
8    D_2 = D;
9    D_2(D(:,C)==2, C) = 1;
10   D_2(D(:,C)~=2, C) = -1;
11
12   D_3 = D;
13   D_3(D(:,C)==3, C) = 1;
14   D_3(D(:,C)~=3, C) = -1;
```

- In last part, we divide the dataset into three different dataset as follow the part (e) requierments, and we modify the $\lambda$ and $T$, compute each $w$ for three dataset and track the binary loss. We select the optimal $w$, which has least binary loss. Then use part (f) function, reclassify the dataset and compute the loss. Repeat this process several times, each time store the optimal weight, binary loss and final loss.
- We obtain the best weight as following:

| w | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| w_1 | -0.9099421 | 2.03256985 | 0.51723937 | 0.58956456 | 1.47053108 | -1.1926746 | -3.9665291 | 0.36748049 |
| w_2 | 8.29929124 | -7.257732 | -1.2886759 | -3.238241 | -1.6755155 | 1.62626611 | -0.3327964 | -1.4819588 |
| w_3 | -4.3032284 | 2.12446638 | 0.34679696 | 1.27768143 | -0.0472385 | 1.70518543 | 2.35108058 | 0.50693703 |

|            | w_1        | w_2        | w_3        |
|------------|------------|------------|------------|
| Binary Loss | 0.09047619 | 0.02857143 | 0.03333333 |

Mutlti-class error: $10/210 = 0.045$