# MACHINE LEARNING
## HOMEWORK - 2
### SHRAINIK JAIN - 132333 B

1) $H(x) = sgn \left\{ \sum_{t=1}^{T} \alpha_t h_t(x) \right\} = sgn\{f(x)\}$

$$f(x) = \sum_{t=1}^{T} \alpha_t h_t(x)$$

$\epsilon_{Training}$ = Probability of wrong prediction over the distribution of weighted data points

$= P_{i \sim D_t} [H_f(x_i) \neq y_i]$

$= \dfrac{\text{Number of training point where we predicted incorrectly}}{\text{Total number of training points}}$

$\leftarrow$ i.e. $1$, when $H(x^i) \neq y^i$
$0$, otherwise.

$\epsilon_{Train} = \dfrac{\sum_{j=1}^{N} \mathbb{1}\{H_o(x^i) \neq y^i\}}{N} \longrightarrow ①$
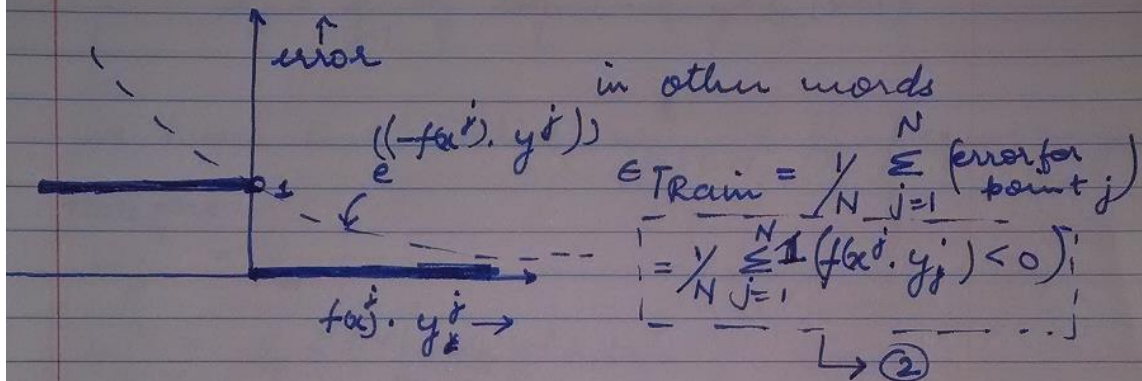
we make an error when
$y^i = 1$ and $f(x^i) \leq 0$ { sgn is -1 }
or
$y^i = -1$ and $f(x^i) \geq 0$ { sgn is +1 }

This implies that we make an error when
$$f(x^i) \cdot y^i \leq 0$$

i.e. error for a single point $= \begin{cases} 1 & , f(x^i) \cdot y_k^i < 0 \\ \\ 0 & , \text{otherwise} \end{cases}$



in other words

$$\epsilon_{Train} = \frac{1}{N} \sum_{j=1}^{N} \left( \begin{array}{c} \text{error for} \\ \text{point } j \end{array} \right)$$

$$= \frac{1}{N} \sum_{j=1}^{N} \mathbb{1}(f(x^i \cdot y_j^i) < 0) \quad \longrightarrow ②$$

The above step function for the error of classifying $j^{th}$ point is bounded by

$\exp(-f(x^i) \cdot y^i)$ because:

$$\exp(-f(x_j) \cdot y^j) > 1 \quad \text{for } -f(x^i) \cdot y^i < 0$$

$$1 \geq \exp(-f(x_j) \cdot y^i) \geq 0 \quad \text{for } -f(x^i) \cdot y^j \geq 0$$

i.e.

$$\exp(-f(x_j) \cdot y^j) \geq f(x^i) \cdot y^i \quad \} \quad \text{for all } (x^i, y^i) \quad \longrightarrow ③$$

combining equations ①, ② & ③ we get:

$$\epsilon_{Training} = \frac{1}{N} \sum_{j=1}^{N} \mathbb{1}\{H(x^i) \neq y^i\} \leq \frac{1}{N} \sum_{j=1}^{N} \exp(-f(x^i) \cdot j)$$

**Q1.** **2)**

$$w_j^{(t+1)} = \frac{w_j^{(t)} \exp(-\alpha_t \, y^i \, h_t(x^i))}{Z_t}$$

$$Z_t = \sum_{j=1}^{N} w_j^{(t)} \exp(-\alpha_t \, y^i \, h_t(x^i))$$

we know that initial $w_{j=1 \text{ to } N}^{1} = 1/N \}$ equal weights

therefore :

$$w_j^{1} = 1/N$$

$$w_j^{2} = 1/N \left( \frac{\exp(-\alpha_1 y^i h_1(x^i))}{Z_1} \right)$$

$$\vdots$$

$$w_j^{(t+1)} = 1/N \left( \frac{\exp(\alpha_1 y^i h_1(x^i))}{Z_1} \right) \left( \frac{\exp(-\alpha_2 y^i h_2(x^i))}{Z_2} \right)$$

$$- \cdots \cdots \cdots \left( \frac{\exp(-\alpha_t y^i h_t(x^i))}{Z_t} \right)$$

$$w_j^{(t+1)} = 1/N \; \frac{\prod_{t=1}^{t} \exp(-\alpha_t y^i h_i(x^i))}{\prod_{i=1}^{t} Z_{ti}}$$

now Sum of all weights at each iteration is 1.

So, $\sum_{j=1}^{N} w_j^{t+1} = 1$

$\therefore \frac{1}{N} \sum_{j=1}^{N} \frac{\left( \prod_{t=1}^{T} \exp(-\alpha_t y^j h_{t,i}(x^j)) \right)}{\prod_{t=1}^{T} z_t} = 1$

$\frac{1}{N} \sum_{j=1}^{N} \exp \left( \sum_{t=1}^{T} (-\alpha_t y^j h_i(x^j)) \right) = \prod_{t=1}^{T} z_t$

$\frac{1}{N} \sum_{j=1}^{N} \exp \left( (-y^j) \left( \sum_{t=1}^{T} \alpha_t h_t(x^j) \right) \right) = \prod_{t=1}^{T} z_t$

$\cdot$ also $\sum_{t=1}^{T} \alpha_t h_t(x^j) = f(x_j)$

$\therefore \boxed{\frac{1}{N} \sum_{j=1}^{N} \exp(-y^j \cdot f(x^j)) = \prod_{t=1}^{T} z_t.}$

Q.E.D.

Q1) 3) (a) $\epsilon_t = \sum\limits_{j=1}^{m} w_j^t \, 1\{h_t(x^j) \neq y^j\}$

$z_t^i = (1-\epsilon_t)\exp(-\alpha_t) + \epsilon_t \exp(\alpha_t)$

$\alpha_t$ for $z_t$ max $\Rightarrow \dfrac{\partial z_t}{\partial \alpha_t} = 0$

$\Rightarrow \quad -(1-\epsilon_t)\exp(-\alpha_t) + \epsilon_t \exp(\alpha_t) = 0$

$\Rightarrow \quad (1-\epsilon_t)\exp(-\alpha_t) = \epsilon_t \exp(\alpha_t)$

$\Rightarrow \quad \ln(1-\epsilon_t) - \alpha_t = \ln(\epsilon_t) + \alpha_t$

$\Rightarrow \quad \ln\left(\dfrac{1-\epsilon_t}{\epsilon_t}\right) = 2\alpha_t$

$$\boxed{\hat{\alpha}_t = \frac{1}{2}\ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)}$$

$z_t^{opt} = (1-\epsilon_t)\exp\left(-\frac{1}{2}\ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)\right)$

$\qquad\qquad + \epsilon_t \exp\left(\frac{1}{2}\ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)\right)$

$\Rightarrow (1-\epsilon_t)\cdot\left(\dfrac{1-\epsilon_t}{\epsilon_t}\right)^{-1/2} + \epsilon_t\left(\dfrac{1-\epsilon_t}{\epsilon_t}\right)^{1/2}$

$\Rightarrow (1-\epsilon_t)^{1/2}\,\epsilon_t^{1/2} + \epsilon_t^{1/2}(1-\epsilon_t)^{1/2}$

$\Rightarrow \boxed{z_t^{opt} = 2\sqrt{\epsilon_t(1-\epsilon_t)}}$ Q.E.D.

Q1 3) (b) $\quad \epsilon_t = \frac{1}{2} - \gamma_t$

$$\Rightarrow z_t^{opt} = 2\sqrt{(\tfrac{1}{2}-\gamma_t)(1-(\tfrac{1}{2}-\gamma_t))}$$

$$\Rightarrow \frac{2}{2}\sqrt{(1-2\gamma_t)(2\gamma_t-1)}$$

$$z_t^{opt} \Rightarrow \sqrt{1-4\gamma_t^2}$$

also, $z_t^{opt} = (1-\epsilon_t)\exp(-\tfrac{\gamma}{2})$

$z_t^{opt} = e^{\frac{1}{2}\ln(1-4\gamma_t^2)}$

also, since $\log(1-x) \le -x$ for $0 \le x \le 1$

$$\ln(1-4\gamma_t^2) \le -4\gamma_t^2$$

because
$0 \le \gamma_t \le \frac{1}{2}$
$\gamma_t^2 \le \frac{1}{4}$
$0 \le 4\gamma_t^2 < 1$

therefore $\Rightarrow$

$$z_t^{opt} \le e^{\frac{1}{2}(-4\gamma_t^2)} = e^{-2\gamma_t^2}$$

$$\therefore z_t\,opt \le \exp(-2\gamma_t^2)$$

$$\therefore z_t \le \exp(-2\gamma_t^2) \quad Q.E.D.$$

$$\therefore \quad \epsilon_{\text{Training}} \leq \prod_{t=1}^{T} z_t \leq \exp\left(-2 \sum_{t=1}^{T} \gamma_t^2\right)$$

Q 1. 3) c) if each classifier is better than random,

then $\quad \epsilon_T < \frac{1}{2}$

$$\epsilon_T = \frac{1}{2} - \gamma_t$$

~~for all~~ $\gamma_t$, for some $\boxed{\gamma_t < \frac{1}{2}}$

this implies $\exists \gamma = \boxed{\min(\gamma_1, \gamma_2 \ldots \gamma_t)}$

$$\therefore \quad \epsilon_{\text{Training}} \leq \exp\left(-2 \sum_{t=1}^{T} \gamma_t^2\right)$$

since $\forall \gamma_t \quad \boxed{\gamma \leq \gamma_t}$

$$\exp(-2\gamma_t^2) \leq \exp(-2\gamma^2)$$

$$\therefore \quad \epsilon_{\text{Train}} \leq \exp(-2\gamma^2) \cdot \exp(-2\gamma^2) \ldots$$
$$\ldots \ldots \exp(-2\gamma^2)$$
$$\underbrace{\qquad\qquad\qquad}_{T \text{ times}}$$

$$\boxed{\epsilon_{\text{Train}} \leq \exp(-2 T\gamma^2)} \quad Q.E.D.$$

## Ques 2. 1

The following figure is the Voronoi diagram for the dataset. The Redline is the decision boundary, i.e. all points left of it are positive.

Q2  1)  Attached as image.

2)  (8,5)
All others change the decision boundary if removed.

3).  Error by removing ~~points~~ points one at a time and classifying it based on other ~~closest~~ data points.

Let the points be numbered 1 to 10 as

positive → | 1 ⇒ (2,2) | 2 ⇒ (2,4) | 3 ⇒ (2,6) | 4 ⇒ (2,8) | 5 ⇒ (4,5)
negative → | 6 ⇒ (4,3) | 7 ⇒ (4,7) | 8 ⇒ (6,4) | 9 ⇒ (6,6) | 10 ⇒ (8,5)

Removing point 1 ⇒.                    ← negative
   Closest points ⇒ (2, 5, 6.) ⇒ prediction
                      ↑  ↑                    positive
                    positive          $\overline{error = 0}$

removing point 2 ⇒   ~~distance~~ distance 2.236
closest points ⇒  (1, 3, 5, 6)     $\overline{2}$
                   ↑  ↑  ↑   ↑
              positive  negative } irrespective of
              prediction          choice of the 3
                                  points
                                    prediction =
                                       positive
                          $\overline{error = 0}$

removing point 3 ⇒  distance = 2
closest points  (2, 4, 5, 7)      → distance = 2.236
                ‾‾‾‾‾‾    ↗
                positive   negative
                prediction = positive
                   $\overline{error = 0}$

removing point 4 ⇒ negative
closest points (3, 5, 7)
      ↑
      positive    prediction = positive

$$\boxed{error = 0}$$

removing point 5 ⇒
                    negative              positive
closest points = ( 6, 7, . 2, 3, 8, 9 )
                    └──┘    └────┬────┘
                  distance=2   distance 2.236

$$\boxed{error = 1}$$

removing point 6 ⇒   distance=2   distance = 2.236
                        ↓         ↗     negative
closest points  = ( 5, 1, 2, 8 }
                    └──┘
                   positive

        prediction = positive

        $$\boxed{error = 1}$$ } as actual
                               value is negative

                    distance 2
removing point 7 ⇒      ↓        distance 2.236
closest points ⇒  ( 5, 3, 4, 9 }
                   └──┬──┘  └ negative
                   positive
    prediction =) positive
        $$\boxed{error = 1}$$     Actual value
                                  = negative

distance 2.236

removing point 8 ⇒ distance 2 positive
closest points ⇒ ( 9, 5, 4, 10 )

negative

prediction = negative
[error = 0].

removing point 9 ⇒ positive
closest points ⇒ ( 8, 5, 7, 10 )

negative

prediction = negative
[error = 0]

removing point 10 ⇒
closest points ⇒ ( 8, 9, 5 )

negative   positive

prediction = negative
[error = 0]

Total error = (1 + 1 + 1) = 3
mean error = 3/10 = 0.3

**Q2. 4)** Removing feature 1:-

Data points ⇒ (2,+) (4,+) (6,+) (8,+) (5,+)
  (3,-) (7,-) (4,-) (6,-) (5,-)

LOOCV iterations

| Point to remove | Closest Points | Output | Error. |
|---|---|---|---|
| (2,+) | (3,-);(4,-),(4,+) | - | 1 |
| (4,+) | (4,-);(3,-);(5,-);(5,+) | - | 1 |
| (6,+) | (6,-);(5,+);(5,-) | - | 1 |
| (8,+) | (7,-);(6,-);(6,+) | - | 1 |
| (5,+) | (5,-);(4,-);(6,-);(4,+);(6,+) | Not Defined | 1 |
| (3,-) | (2,+);(4,+);(4,-) | + | 1 |
| (7,-) | (6,+);(6,-);(8,+) | + | 1 |
| (4,-) | (4,+);(3,-);(5,-);(5,+) | Not Defined | 1 |
| (6,-) | (6,+);(5,+);(5,-) | + | 1 |
| (5,-) | (5,+);(4,-);(4,+);(6,-);(6,+) | Not Defined | 1 |

Total error = 10

Removing feature 2:

Data points ⇒ (2,+) (2,+) (2,+) (2,+) (4,+) (4,-) (4,-) ,(6,-) ,(6,-)
  (8,-)

LOOCV Iterations.

| Point to remove | Closest points | Output | Error |
|---|---|---|---|
| (2,+) | (2,+);(2,+);(2,+) | + | 0 |
| Same for other points at (2,+) | | + | 0 |
| (4,+) | (4,-);(4,-);[(2,+)] ×4;(6,-);(6,-) | Not defined | 1 |

    4 points with (2,+)

| | | | |
|---|---|---|---|
| (4,-) | (4,+),(4,-)[(2,+)]×4; (6,-);(6,-) | Not defined | 1 |
| (4,-) | ← Same → | Not defined | 1 |
| (6,-) | (6,-);(4,-);(4,-);(4,+);(8,-); | - | 0 |
| (6,-) | ← Same → | - | 0 |
| (8,-) | (6,-);(6,-);(4,-);(4,-);(4,+) | - | 0 |

Total ⇒ 3

THEREFORE, we can safely eliminate feature 2 .

Q3. 1) Let Y = has college degree and N = Does not have college degree

Root ⇒ | Y=6 |
        | N=4 |

$H(Y) = -\frac{6}{10} \log \frac{6}{10} - \frac{4}{10} \log \frac{4}{10}$

$\Rightarrow 0.970951.$

Step 2:

we have the following intuitive splits
① salary ≤ 27000
② salary ≥ 65000
③ age ≥ 43

$H[Y|X]$

IG for choice 1 ⇒ $0.9709 - \left[ -0.8 \left( \frac{2}{8} \log \frac{2}{8} \right) + \frac{6}{8} \log (6/8) - 0.2 \times 0 \right]$

⇒ 0.322

$H[Y|X]$

IG for choice 2 ⇒ $0.9709 - \left[ -0.3 \times 0 - 0.7 \left( \frac{3}{7} \log \frac{3}{7} + \frac{4}{7} \log (4/7) \right) \right]$

⇒ 0.282

IG for choice 3 ⇒ $0.9709 - \left[ -0.5 \left( \frac{2}{5} \log \frac{2}{5} + \frac{3}{5} \log \frac{3}{5} + \frac{1}{5} \log \frac{1}{5} + \frac{4}{5} \log \frac{4}{5} \right) \right]$

⇒ 0.125560.

Max information gain for choice 1:

Decision
stump
after step 1:

| Y=6 |
| N=4 |                          | IG = 0.322 |

salary > 27000          salary ≤ 27000
P = 0.8                 P = 0.2

| Y=6 |                 | N=2 |
| N=2 |                 | Y=0 |

                        predict "N"

Step 2:
Now we have the following intuitive splits ⇒
① age ≥ 43
② salary ≥ 65000

$\begin{array}{|c|}\hline Y=6\\ N=4\\\hline\end{array}$

salary $\geq 27000$       salary $\leq 27000$

$P = 0.8$            $P = 0.2$

$\begin{array}{|c|}\hline Y=6\\ N=2\\\hline\end{array}$        $\begin{array}{|c|}\hline Y=0\\ N=2\\\hline\end{array}$

           Predict "N"
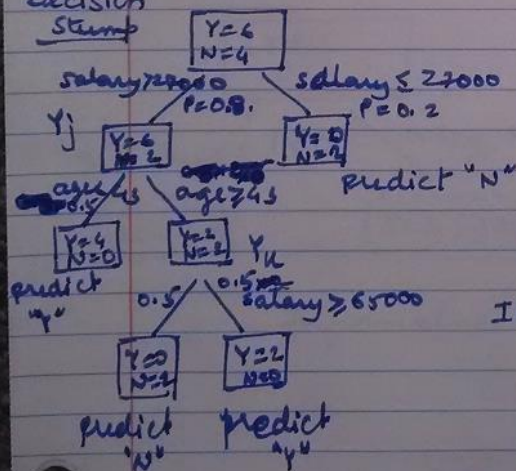
$H(Y \mid \text{salary} \geq 27000) = 0.649022.$

$H(Y_j) \Rightarrow -1\left(\frac{2}{8}\log\frac{2}{8} + \frac{6}{8}\log\frac{6}{8}\right) = 6.811)$

# IG for choice1 $\Rightarrow H(Y_j) - \left[0.5\times 0 + 6.5(\frac{1}{2}\log_2 0.5 + 0.5\log_2 0.5)\right]$
         $\Rightarrow H(Y_j) - 0.5$
           $\Rightarrow 0.311$

IG for choice 2 $\Rightarrow H(Y_j) - \left[-\frac{3}{8}\times 0 + -\frac{5}{8}\left(\frac{2}{5}\log(2/5) + 3/5\log(3/5)\right)\right]$
         $\Rightarrow H(Y_j) - 0.6068$
           $\Rightarrow 0.2042.$

Choice 1 has higher information gain $\Rightarrow$ $\boxed{0.311}$

Decision Stump

$\begin{array}{|c|}\hline Y=6\\ N=4\\\hline\end{array}$

salary $\geq 27000$    salary $\leq 27000$

$P = 0.8.$       $P = 0.2$

$Y_j$   $\begin{array}{|c|}\hline Y=6\\ N=2\\\hline\end{array}$    $\begin{array}{|c|}\hline Y=0\\ N=2\\\hline\end{array}$

age$\leq 45$   age $\geq 45$   Predict "N"

$\begin{array}{|c|}\hline Y=4\\ N=0\\\hline\end{array}$   $\begin{array}{|c|}\hline Y=2\\ N=2\\\hline\end{array}$   $Y_u$

predict
"y"   0.5 / 0.5   salary $\geq 65000$

$\begin{array}{|c|}\hline Y=0\\ N=2\\\hline\end{array}$   $\begin{array}{|c|}\hline Y=2\\ N=0\\\hline\end{array}$

predict   predict
"N"     "y"

**Step 3:**

After step 2: there is just one intuitive step choice
$\rightarrow$ salary $\geq 65000$

$H(Y_k) = -1\times\left(\frac{2}{4}\log\frac{2}{4} + 2/4\log 2/4\right)$
       $\Rightarrow 1$

IG for the only choice
$\Rightarrow H(Y_k) - \left[\frac{1}{2}\times 0 + \frac{1}{2}\times 0\right]$
$\Rightarrow (1 - 0) . \Rightarrow \boxed{1}$

$\boxed{\text{Depth of Tree} = 3}$

**Q3.2)**

Q3.2) decision based on $\alpha x_{age} + \beta x_{income} - 1$

Information gain is maximum when feature perfectly separated.

writing the decision boundary based on points as the equation of line in the form

$$y = mx + c$$

$$x_{income} = \left(-\frac{\alpha}{\beta}\right) x_{age} + \frac{1}{\beta}$$

Solving for $\alpha$ and $\beta$ by putting in points
(24, 40000) and (22, 32000)

we get

$$\boxed{\alpha = -\frac{1}{16} \qquad \beta = \frac{1}{16000}}$$



$$Y = 6$$
$$N = 4$$

$\alpha x_{age} + \beta x_{income} - 1 < 0$   $\qquad \alpha x_{age} + \beta x_{income} - 1 \geqslant 0$   $\left.\right\}$ $\boxed{\text{Depth of tree} = 1}$

$\boxed{N=4 \quad Y=0}$  $\qquad$  $\boxed{Y=6 \; ; \; N=0}$

Information gain $= H(Y) - H(Y|x)$
$$\Rightarrow 0.970951 - 0$$
$$\boxed{IG \Rightarrow 0.97095}$$

**Ques 3.3**

Advantages of multivariate decision trees:
- More than one feature can be tested for per decision, this helps getting a way better predictor when the data points are linearly separable.
- Training error is lower as more complex models are allowed (decision based on multiple variables and rather than a single variable).
- Multivariate decision trees will have, on an average, smaller tree size.

Disadvantages of multivariate decision trees:
- The computation cost (CPU Time) of the learning algorithm is very high (imagine running linear regression again and again for each decision node), as a lot of features need to be tested to find out the maximum information gain per step. While the univariate algorithm requires considering only one feature at a time.

- On smaller datasets, multivariate trees tend to over-fit a lot because of standard deviation in the data points. This leads to higher test errors. Multivariate trees require a lot of data-points which might or might not be present.
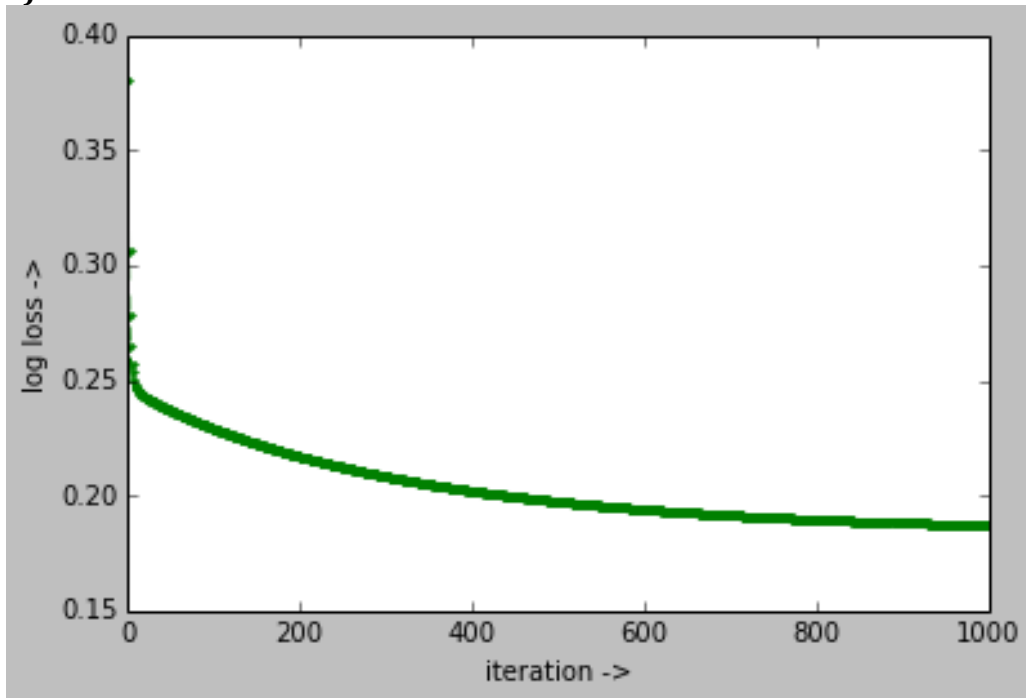
**Ques 4.4.1**

Actual weight update step:

$$w_i^{t+1} = w_i^t + \eta \left\{ -\lambda w_i^t + \frac{1}{N} \left( \sum_{j=1}^{N} x_i^j \left[ y^j - \frac{e^{w_0 + \Sigma w_i x_i^j}}{1 - e^{w_0 + \Sigma w_i x_i^j}} \right] \right) \right\}$$

Weight update step as in python (Y is Nx1, X is Nx(D+1), W is (D+1)x1):

yj_minus_p = Y – (exp(X.dot(W))/(1+ exp(X.dot(W))))
W[0] = W[0] + (eta / N) * sum(yj_minus_p)
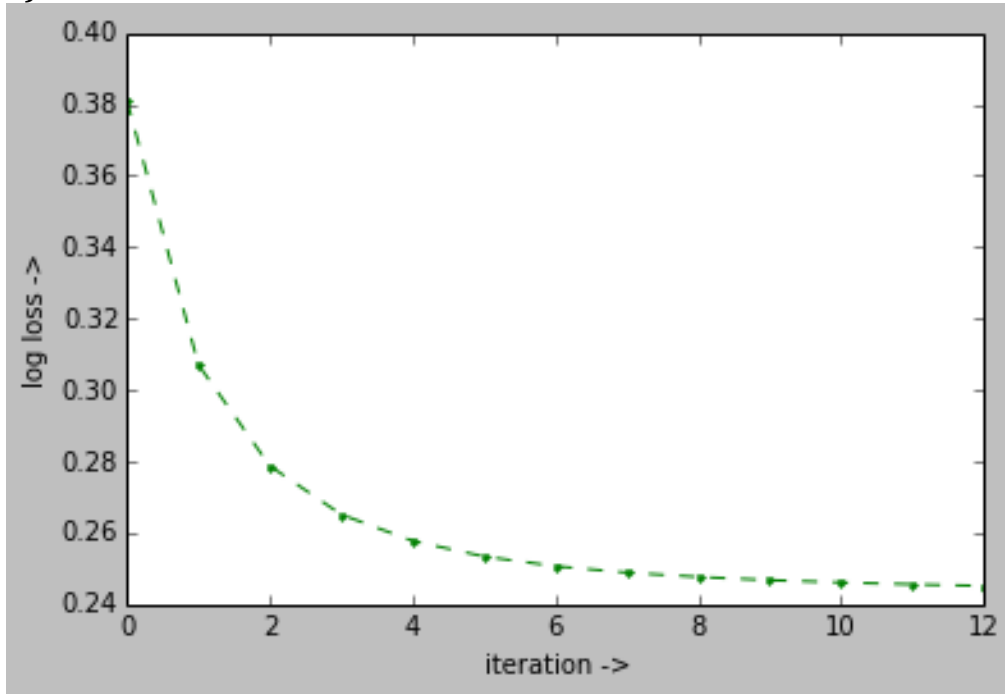W[1:] = (1 - eta * lmbd) * W[1:] + eta * (array(yj_minus_p).dot(X[:,1:]) / N)

**Ques 4.4.2**
**a)**



**b)** SSE for batch gradient descent with 1000 iterations: 54.0

**Ques 4.4.3 a)** Number of iterations with the stopping criteria: 13

**b)**



**c)** SSE for batch gradient descent with the stop criteria: 54.0

**Ques 4.5.1**

Actual weight update step:

$$w_i^{t+1} = w_i^t + \eta\left\{-\lambda w_i^t + x_i^j\left[y^j - \frac{e^{w_0 + \Sigma w_i x_i^j}}{1 - e^{w_0 + \Sigma w_i x_i^j}}\right]\right\}$$

Weight update step as in python (Y is Nx1, X is Nx(D+1), W is (D+1)x1):

yj_minus_p = Y[j] - (1-1/(1+exp(X[j,:].dot(W))))
W[o] = W[o] + (eta) * yj_minus_p
W[1:] = (1 - eta * lmbd) * W[1:] + eta * (yj_minus_p * (X[j,1:]))

**Ques 4.5.2**
**a)** $L_2$ Norm for lambda = 0 is: 1.92503456434
$L_2$ Norm for lambda = 0.3 is: 0.283520413611

**b)** SSE for lambda = 0.3 is: 54.0

**c)** Feature Weights for INTERCEPT: -3.10616785425
DEPTH: 0.109353101677
POSITION: -0.006094751226

**Ques 4.5.3**
After 5 iterations, log loss with:
Stochastic Descent: 0.197392978738
Gradient Descent: 0.257743788383

Stochastic Descent seems to converge faster.

**Ques 4.6.1**
For predictions made by SGD running with one pass over data and lambda = 0.3 and eta = 0.1:
Precision and Recall for class 0:      0.946,          1.0
Precision and Recall for class 1:      0,              0.0

**Ques 4.6.2**
For predictions made by batch gradient descent running for 10000 iterations over the oversampled data and with lambda = 0.3 and eta = 0.01:

Precision and Recall for class 0:      0.94140625          0.509513742072
Precision and Recall for class 1:      0.049180327         0.444444444444