

CSE 546 Machine Learning, Autumn 2013
Homework 2
Shrainik Jain - 1323338

1 Boosting

MACHINE LEARNING

HOMEWORK - 2

SHRAINK JAIN - 132 3338

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

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

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

$$= P_{i \in D_t} [H_t(x_i) \neq y_i]$$

= Number of training point where we predicted incorrectly

Total number of training points

$$\epsilon_{\text{Train}} = \frac{\sum_{j=1}^N \mathbb{1}_{\{H_0(x_j) \neq y_j\}}}{N} \quad \begin{array}{l} \leftarrow \text{ie. } 1, \text{when } H_0(x_j) \neq y_j \\ 0, \text{otherwise.} \end{array} \rightarrow ①$$

we make an error when

$$y^i = 1 \text{ and } f(x^i) \leq 0 \quad \{ \operatorname{sgn} \text{ is } -1 \}$$

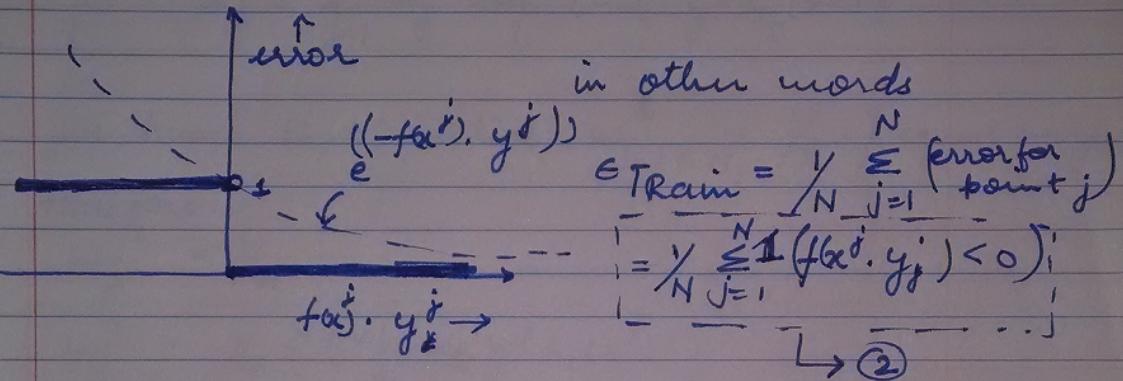
$$\text{or} \quad y^i = -1 \text{ and } f(x^i) \geq 0 \quad \{ \operatorname{sgn} \text{ is } +1 \}$$

This implies that we make an error

when

$$|f(x^i) \cdot y^i| \leq 0$$

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



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

$\exp(-f(x^j) \cdot y_j^j)$ because:

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

$$1 \geq \exp(-f(x^j) \cdot y_j^j) \stackrel{\text{as}}{\geq} 0 \quad \text{for } -f(x^j) \cdot y_j^j \geq 0$$

i.e.

$$\exp(-f(x^j) \cdot y_j^j) \geq f(x^j) \cdot y_j^j \quad \left\{ \begin{array}{l} \text{for all} \\ (x^j, y_j^j) \end{array} \right\}$$

combining equations ①, ② & ③ we get:

$$e_{\text{Training}} = \frac{1}{N} \sum_{j=1}^N \mathbb{1}_{\{H(x^j) \neq y_j^j\}} \leq \frac{1}{N} \sum_{j=1}^N \exp(-f(x^j) \cdot y_j^j)$$

$$\text{Q 1. 2)} \quad w_j^{(t+1)} = \frac{w_j^{(t)}}{\sum_{j=1}^N w_j^{(t)} \exp(-\alpha_t y^j h_t(x^j))}$$

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

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

therefore:

$$w_j^1 = \frac{1}{N}$$

$$w_j^2 = \frac{1}{N} \left(\frac{\exp(-\alpha_1 y^j h_1(x^j))}{z_1} \right)$$

$$w_j^{(t+1)} = \frac{1}{N} \left(\frac{\exp(-\alpha_1 y^j h_1(x^j))}{z_1} \cdot \frac{\exp(-\alpha_2 y^j h_2(x^j))}{z_2} \cdot \dots \cdot \frac{\exp(-\alpha_t y^j h_t(x^j))}{z_t} \right)$$

$$w_j^{(t+1)} = \frac{1}{N} \frac{\prod_{i=1}^t \exp(-\alpha_i y^j h_i(x^j))}{\prod_{i=1}^t z_{bi}}$$

now sum of all weights at each iteration
is 1.

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

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

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

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

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

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

Q.E.D.

Q1) 3) (a) $\epsilon_t = \sum_{j=1}^m w_j^{t-1} \{h_t(x_j) \neq y_j\}$

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

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

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

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

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

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

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

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

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

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

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

$$\begin{aligned}
 Q1 3) \quad b) \quad \epsilon_t &= \frac{1-y_t}{2} \\
 \Rightarrow z_t^{\text{opt}} &= 2 \sqrt{(y_2 - y_t)(1 - (y_2 - y_t))} \\
 &= \frac{2}{2} \sqrt{(1-2y_t)(2y_t-1)} \\
 z_t^{\text{opt}} &\equiv \sqrt{1-4y_t^2}
 \end{aligned}$$

also, $z_t^{\text{opt}} = (1-\epsilon_t) \exp(-\gamma_t)$

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

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

$$\ln(1-4y_t^2) \leq -4y_t^2$$

because
 $0 \leq y_t \leq \frac{1}{2}$

$$y_t^2 \leq \frac{1}{4}$$

$$0 \leq -4y_t^2 < 1$$

therefore $\Rightarrow z_t^{\text{opt}} \leq e^{\frac{1}{2}(-4y_t^2)} = e^{-2y_t^2}$

$$\therefore z_t^{\text{opt}} \leq \exp(-2y_t^2)$$

$$\therefore [z_t] \leq \exp(-2y_t^2) \quad \text{Q.E.D.}$$

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

Q 1. 3) c) if each classifier is better than random,
then $\epsilon_T < \frac{1}{2}$

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

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

$$\text{this implies } \exists \gamma = \min(\gamma_1, \gamma_2, \dots, \gamma_T)$$

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

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

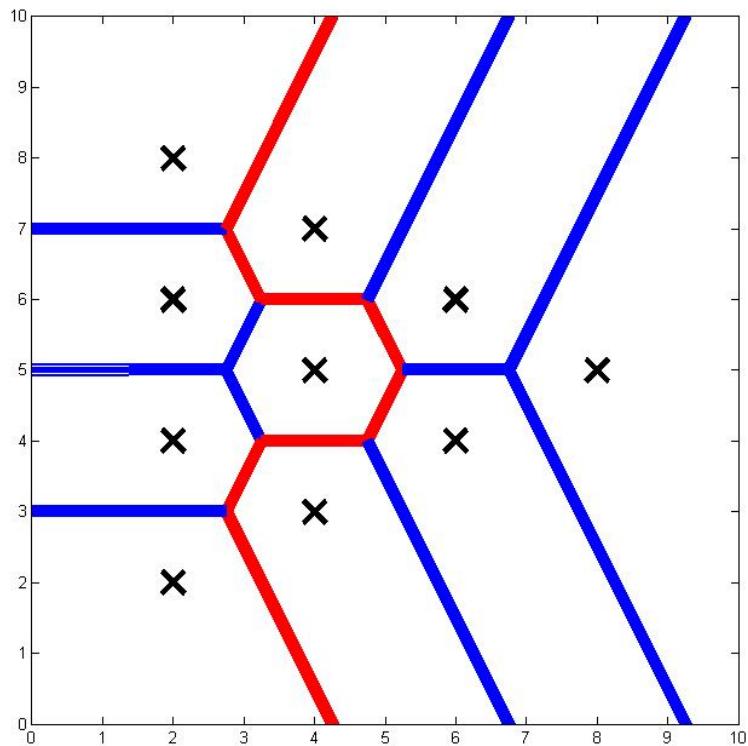
$$\exp(-2 \gamma_t^2) \leq \exp(-2 \gamma^2)$$

$$\therefore \epsilon_{\text{Training}} \leq \underbrace{\exp(-2 \gamma^2)}_{\text{exp}(-2 \gamma^2) \dots \exp(-2 \gamma^2)} \dots$$

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

2 KNN

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 ~~(8, 5)~~ points one at a time and classifying it based on other ~~others~~ data points.

Let the points be numbered 1 to 10 as
positive \Rightarrow 1 \Rightarrow (2, 2) | 2 \Rightarrow (2, 4) | 3 \Rightarrow (2, 6) | 4 \Rightarrow (2, 8) | 5 \Rightarrow (4, 5)
negative \Rightarrow 6 \Rightarrow (4, 3) | 7 \Rightarrow (4, 4) | 8 \Rightarrow (6, 4) | 9 \Rightarrow (6, 6) | 10 \Rightarrow (8, 5)

removing point 1 \Rightarrow negative

closest points \Rightarrow (2, 5, 6) \Rightarrow prediction
position position $\boxed{\text{Error} = 0}$

removing point 2 \Rightarrow distance $= \sqrt{2}$ distance $= 2.236$

closest points \Rightarrow (1, 3, 5, 7)

position negative } inspection of
prediction choice of the 3
points
 $\boxed{\text{Error} = 0}$

removing point 3 \Rightarrow distance $= \sqrt{2}$ distance $= 2.236$

closest points \Rightarrow (2, 4, 5, 7)

position negative
prediction = positive
 $\boxed{\text{Error} = 0}$

removing point 4 \Rightarrow negative
 closest points $(3, 5, 7)$
 positive prediction = positive
 $\boxed{\text{error} = 0}$

removing point 5 \Rightarrow
 closest points = $(6, 7, 2, 3, 8, 9)$
 distance = 2 distance 2.236
 $\boxed{\text{error} = 1}$

removing point 6 \Rightarrow distance 2 distance 2.236
 closest points = $(5, 1, 2, 8)$
 position negative
 prediction = positive
 $\boxed{\text{error} = 1}$ } as actual value is negative

removing point 7 \Rightarrow distance 2 distance 2.236
 closest points $\Rightarrow (5, 3, 4, 9)$
 position negative
 prediction = positive
 $\boxed{\text{error} = 1}$ Actual value = negative

distance 2.236

removing point 8 \Rightarrow distance 2 positive
 closest points $\Rightarrow (9, 5, 6, 10)$

negative
 prediction = negative
 $\boxed{\text{error} = 0}$.

removing point 9 \Rightarrow
 closest points $\Rightarrow (8, 5, 7, 10)$

positive
 negative
 prediction = negative
 $\boxed{\text{error} = 0}$

removing point 10 \Rightarrow
 closest points $\Rightarrow (8, 9, 5)$

negative positive
 prediction = negative
 $\boxed{\text{error} = 0}$

$$\begin{aligned}\text{Total error} &= (1 + 1 + 1) = 3 \\ \text{mean error} &= 3/10 = 0.3\end{aligned}$$

Q2. 4)

Removing feature 1:-

Data points $\Rightarrow (2, +) (4, +) (6, +) (8, +) (5, +)$

$(3, -) (7, -) (4, -) (6, -) (5, -)$

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 $\Rightarrow (2, +) (2, +) (2, +) (2, +) (4, +) (4, -) (4, -) (6, -) (6, -)$
 $(8, -)$

LOOCV Iterations.

Point to remove	Closest points	Output	Error
$(2, +)$	$(2, +); (2, +); (2, +)$	+	0
Same for other points at $(2, +)$		+	0
$(4, +)$	$(4, -); (4, -); \underbrace{(2, +)}_{\text{4 points with } (2, +)}; (4, +); (6, -); (6, -)$	Not defined	1
$(4, -)$	$(4, +); (4, -); [(2, +) \times 4]; (6, -); (6, -)$	Not defined	1
$(4, -)$	$\leftarrow \text{Same} \rightarrow$	Not defined	1
$(6, -)$	$(6, -); (4, -); (4, -); (4, +); (8, -);$	-	0
$(6, -)$	$\leftarrow \text{Same} \rightarrow$	-	0
$(8, -)$	$(6, -); (6, -); (4, -); (4, -); (4, +)$	-	0
			Total $\Rightarrow 3$

THEREFORE, we can safely eliminate feature 2.]

3 Decision Trees

3.1

Q 3.1) Let $Y = \text{has college degree}$ and $N = \text{does not have college degree}$

$$\text{Root} \Rightarrow \begin{bmatrix} Y=6 \\ N=4 \end{bmatrix} \quad H(Y) = -\frac{6}{10} \log \frac{6}{10} - \frac{4}{10} \log \frac{4}{10}$$

$$\Rightarrow 0.970951.$$

Step 1:

We have the following intuitive splits

- ① salary ≤ 27000
- ② salary > 65000
- ③ age ≥ 43

$$H[Y|X]$$

$$\text{IG for choice 1} \Rightarrow 0.9709 - [-0.8(\frac{2}{8} \log \frac{2}{8}) + \frac{6}{8} \log(\frac{6}{8})] = 0.2 \times 0 \\ \Rightarrow 0.322$$

$$H[N|X]$$

$$\text{IG for choice 2} \Rightarrow 0.9709 - [-0.3 \times 0 - 0.7(\frac{3}{7} \log \frac{3}{7} + \frac{4}{7} \log \frac{4}{7})] \\ \Rightarrow 0.282$$

$$\text{IG for choice 3} \Rightarrow 0.9709 - [0.5(\frac{4}{5} \log \frac{4}{5} + \frac{3}{5} \log \frac{3}{5} + \frac{1}{5} \log \frac{1}{5} + \frac{1}{5} \log \frac{1}{5})] \\ \Rightarrow 0.125560.$$

Max information gain for choice 1:

Decision

Stump
after Step 1:

$$\text{salary} > 27000 \quad P=0.8 \quad \text{salary} \leq 27000 \quad P=0.2$$

$$\begin{bmatrix} Y=6 \\ N=2 \end{bmatrix}$$

$$\text{IG} = 0.322$$

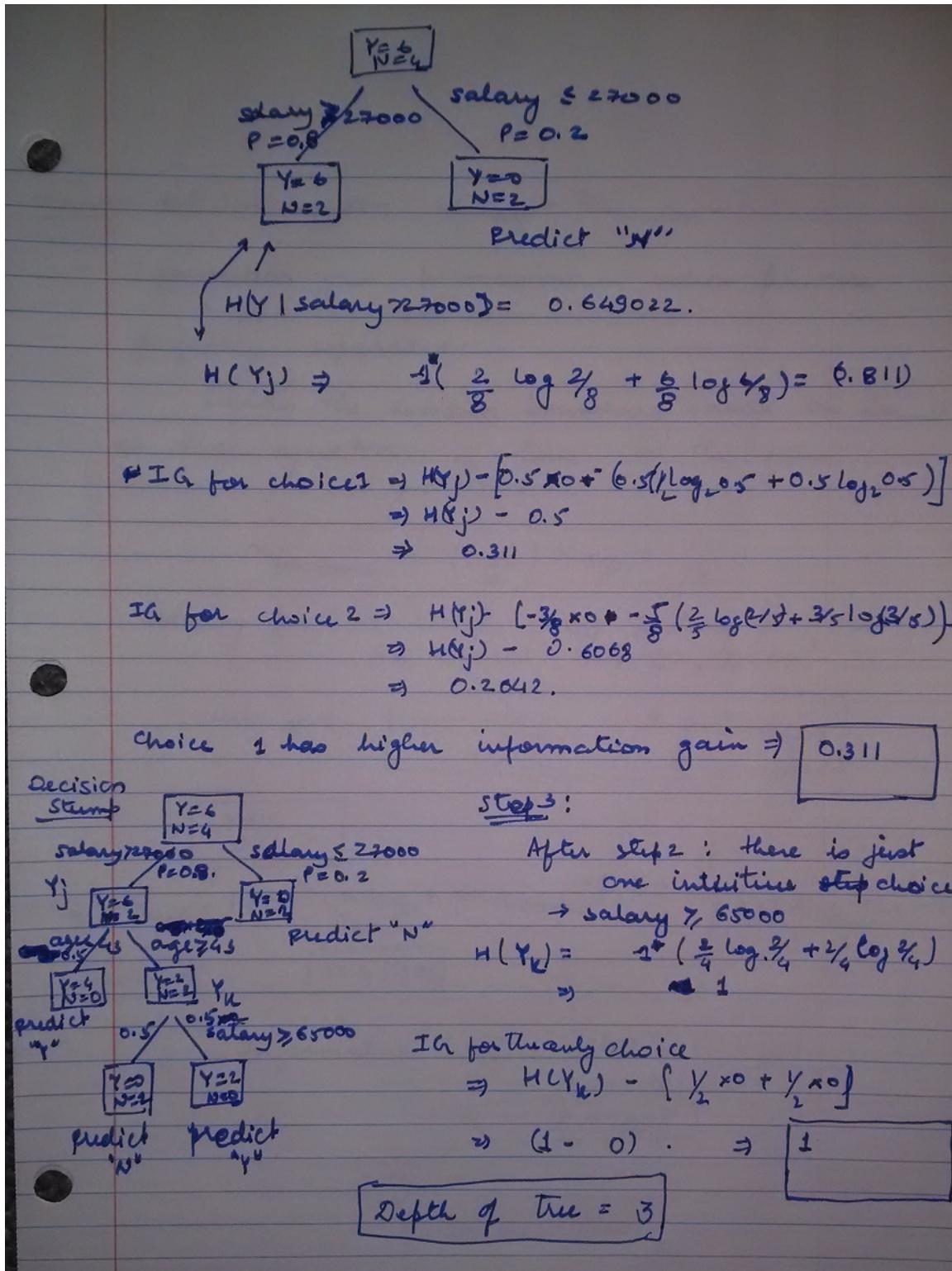
$$\begin{bmatrix} Y=6 \\ N=0 \end{bmatrix}$$

Step 2:

Now we have the following intuitive splits \Rightarrow

- ① age ≥ 43
- ② salary ≥ 65000

predict "N"



3.2

~~Q3.1)~~ ^(Q3.2) decision based on $\alpha \text{age} + \beta \text{income} - 1$

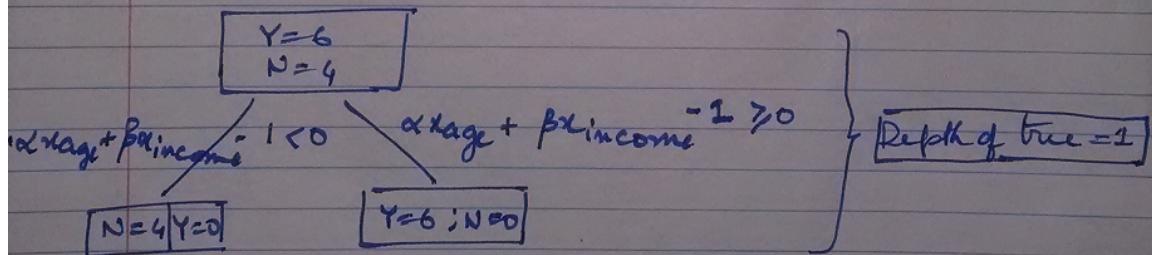
Information gain is maximum when points perfectly separated.

Writing the decision boundary based on points as the equation of line in the form $y = mx + c$

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

Solving for α and β by putting in points (24, 40000) and (22, 38000)

we get $\alpha = -1/16 \quad \beta = 1/16000$



$$\text{Information Gain} = H(Y) - H(Y|x)$$

$$\Rightarrow 0.970951 - 0$$

$$16 \Rightarrow 0.97095$$

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.

4 Programming Question

4.1 Unprocessed Dataset

No question in this part.

4.2 Processed Dataset (that you will use...)

No question in this part.

4.3 Accessing and Processing the Data

No question in this part.

4.4 Batch Gradient Descent

1. Actual weight update step:

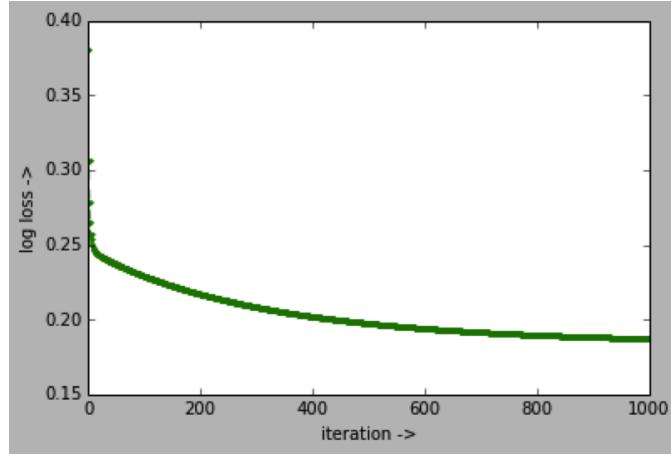
$$w_i^{t+1} = w_i^t + \eta(-\lambda w_i^t + \sum_{j=1}^N x_i^j (y^j - \frac{e^{w_0 + \sum w_i x_i^j}}{1 + e^{w_0 + \sum w_i x_i^j}})) \quad (1)$$

Weight update step in python:

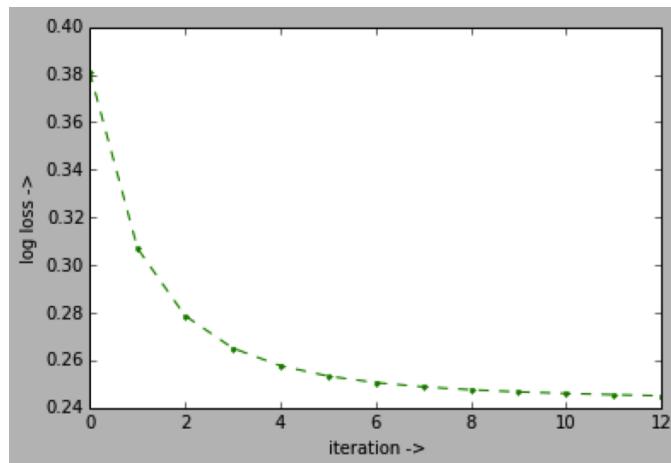
```
#(Y is Nx1, X is Nx(D+1), W is (D+1)x1)

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)
```

2. (a) Log loss versus number of iterations:



- (b) SSE for batch gradient descent with 1000 iterations: 54.0
3. (a) Number of iterations with the stopping criteria: 13
 - (b) Log loss versus Iterations for gradient descent with the stop criteria:



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

4.5 Stochastic Gradient Descent

1. Actual weight update step:

$$w_i^{t+1} = w_i^t + \eta(-\lambda w_i^t + x_i^j(y^j - \frac{e^{w_0 + \sum w_i x_i^j}}{1 + e^{w_0 + \sum w_i x_i^j}})) \quad (2)$$

Weight update step 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[0] = W[0] + (eta) * yj_minus_p
W[1:] = (1 - eta * lmbd) * W[1:] + eta * (yj_minus_p * (x[j,1:]))

```

2. (a) l_2 norm for:

- $\lambda = 0 : 1.92503456434$
- $\lambda = 0.3 : 0.283520413611$

(b) SSE for $\lambda = 0.3 : 54.0$

(c) Feature weights for:

- *INTERCEPT* : -3.10616785425
- *DEPTH*: 0.109353101677
- *POSITION*: -0.006094751226

3. After 5 iterations, log loss with:

- Stochastic Descent: 0.197392978738
- Gradient Descent: 0.257743788383

Stochastic Descent seems to converge faster.

4.6 Class Imbalance

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, 0.0
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, 0.509513742
 - Precision and Recall for class 1: 0.04918032, 0.444444444