

Theory Questions

Yair Shoshan

① $f^{(t+1)}(z_i) = z_{i+1} = h(w^{(t+1)} \cdot z_i + b^{(t)})$

 $h(z_i) = \text{Sign}(z_i)$
 $f \in \mathcal{F}$

(a) we saw that for H which is the class of all linear separators in \mathbb{R}^d :

I $\text{VCdim}(H) = d+1$

so from Sauer-Shelah we have for $n \geq d+1$

$$\mathcal{P}_H(n) \leq \left(\frac{en}{d+1}\right)^{d+1}$$

II we also proved (Ex. 4) that for concatenation $F = F_1 \times F_2$

$$\mathcal{P}_F(n) \leq \mathcal{P}_{F_1}(n) \cdot \mathcal{P}_{F_2}(n)$$

generalizing it to f which is a cartesian product of linear separators

$$f \in \{l_1 \times l_2 \times \dots \times l_d \mid l_i \in H\}$$

by induction! Define $P^{(1)} = \{l_1 \times l_2 \times \dots \times l_k \mid l_i \in H\}$
then the induction assumption:

$$\mathcal{P}_{P^{(k)}} \leq (\mathcal{P}_H(n))^k$$

now define $P^{(k+1)} = f^{(k)} \times l^{(k+1)}$ $l^{(k+1)} \in H$

$$\Rightarrow \Pi_{F^{k+1}} \leq \Pi_F^k \cdot \Pi_F \leq (\Pi_H(n))^d \Pi_H(n) \\ = \varrho \Pi_H(n)^{d+1}$$

\Rightarrow we have for t :

$$\Pi_F(n) \leq (\Pi_H(n))^d \leq \left(\frac{\varrho n}{d+1}\right)^{(d+1)d}$$

b) we saw that for composition:

$$F = F_1 \circ F_2 \Rightarrow$$

$\Pi_{F(n)} \leq \Pi_{F_1}(n) \Pi_{F_2}(n)$. so if we have a function class C :

$C = H \circ F_1 \circ F_2 \dots F_{L-1}$ we can generalize the lemma above (easily as before, by induction)

$$\begin{aligned} \Pi_C &= \Pi_{H \circ F_1 \circ F_2 \dots F_{L-1}} \leq \Pi_H(n) \cdot \Pi_{F_1}(n) \dots \Pi_{F_{L-1}}(n) \\ &\leq \left(\frac{\varrho n}{d+1}\right)^{d+1} \left(\frac{\varrho n}{d+1}\right)^{(d+1)d(L-1)} = \left(\frac{\varrho n}{d+1}\right)^{(d+1)(dL-d+1)} \end{aligned}$$

c) each hidden layer is transformed to the next hidden layer by a function t which is composed defined as the sign of an input calculated with a matrix $W \in \mathbb{R}^{d \times d}$ and a vector $b \in \mathbb{R}^d$ so that $t: \mathbb{R}^d \rightarrow \mathbb{R}^d \Rightarrow d^2+d$ for each hidden layer. the last layer though,

is calculated by a function

$f^L : \mathbb{R}^d \rightarrow \mathbb{R}$ and for that we need
a vector $w^L \in \mathbb{R}^d$ and a bias number
 $b^L \in \mathbb{R}$

$f^{L+1} = L(w^L z_L + b^L)$ so we actually
need $d+1$ parameters for that layer

we thus have in total

$$(L-1)(d^2 + d) + d + 1 = (L-1)d^2 + Ld + 1$$

all hidden layers last layer parameters

(d) first we notice that for two events
 A, B .

$$A \Rightarrow B \text{ iff } \bar{B} \Rightarrow \bar{A}$$

$$\text{so } (2^n \leq (\epsilon n)^n \Rightarrow 2^n \log_2(\epsilon n) \geq n) \\ \Leftrightarrow (n > 2^n \log_2(\epsilon n) \Rightarrow 2^n > (\epsilon n)^n)$$

Let's begin: $2^n > (\epsilon n)^n$

$$n > n \log_2(\epsilon n)$$

$$\frac{n}{n \log_2(\epsilon n)} > 1$$

$$\text{define } f(x) = \frac{x}{\log_2(e x)}$$

$$\frac{dt}{dx} = \frac{1}{N} \left(\log_2(e^x) - \log_2(e) \right) \Rightarrow$$

$$\frac{dt}{dx} \Big|_{x>1} > 0$$

$$\Rightarrow \text{for } x_2 > x_1 \Rightarrow f(x_2) > f(x_1)$$

$$\text{so we can set } x_2 = n$$

$$x_1 = 2N \log_2(eN)$$

and because we start from $x_2 > x_1$
we have $f(x_2) > f(x_1) \Rightarrow$

$$\frac{n}{n \log_2(eN)} > \frac{2N \log_2(eN)}{\log_2(eN \log_2(eN))} =$$

$$\frac{2 \log_2 eN}{1 + \log_2(eN) + \log_2(\log_2(eN))} > 1$$

* we shall prove

$$\log_2(eN) > 1 + \log_2(\log_2(eN))$$

$$\Rightarrow eN > 2 \log_2 eN$$

$$e > \log_2(2e) \quad (\text{True for } N=2)$$

$$e - \log_2(2e) > 0$$

$$\text{define } a(x) = x - 2 \log_2(2x)$$

$$g(x) = x - \frac{2}{x \ln(2)}$$

$$g'(x \geq 2) > 0$$

\Rightarrow for $n \geq 2$:

$$en > 2 \log_2 en$$

$$\Rightarrow \frac{n}{\log_2(en)} > 1$$

$$n > \log_2(en)^n \Rightarrow$$

$$2^n > (en)^n \quad (\text{for } n \geq 2)$$

(for $d=1$ $2^n > en$ for $n > 3$)

\Rightarrow for $n \geq 1$ or $n = 1$ $n \geq 3$

we proved that

$$2^n \leq (en)^n \rightarrow 2^n \log_2(en) \geq n$$

② $\mathcal{N}_c(n) \leq \left(\frac{en}{d+1}\right)^{(d+1)(dL-d+1)}$ (from b)

From c we have $n = (L-1)d^2 + Ld + 1$
 $= (d+1)(dL-d+1)$

$$\Rightarrow \mathcal{N}_c(n) \leq \left(\frac{en}{d+1}\right)^n \leq (en)^n$$

\uparrow
 $d > 0$

$$Vcdim(C) \geq Vcdim(H) = d+1$$

(HCC)

all linear classifiers could be implemented with $L=1$

$$\Rightarrow 2^{\text{VCdim}(C)} = \mathcal{P}_c(\text{VCdim}(C)) \leq (e \cdot \text{VCdim}(C))^N$$

$$\xrightarrow[\text{from d}]{\quad} \text{VCdim}(C) \leq 2N \log_2(eN)$$

	x_1	x_2	x_3	
(1)	Data	(1 1 1)	1	
		(1 0 0)	1	
		(1 1 0)	0	
		(0 0 1)	0	

(a) According to the ID3 Algorithm
each split j will maximize $G(S, i)$

$$j = \arg \max_{i \in A} G(S, i) \quad A \text{ is all possible splits.}$$

$\sim G(S, i) = H(Y) - H(Y|X_i)$ and
 $H(Y) = -\sum_i p(y_i) \log(p(y_i))$

Before the first split we have $P(y=1) = P(y=0) = \frac{1}{2}$
 $\Rightarrow H(Y) = (-\frac{1}{2} \log \frac{1}{2}) \cdot 2 = 1$

now we shall calculate $H(Y|X_i)$ for all i 's

$$H(Y|X_1) = H(Y|X_1=1)P(X_1=1) + H(Y|X_1=0)P(X_1=0)$$

$$\stackrel{!}{=} \frac{3}{5} \left(-\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} \right) + \frac{1}{5} (-1 \log 1) \approx 0.69$$

$$H(Y|X_2) = \frac{1}{2} \left(-\frac{1}{2} \log \frac{1}{2} \right) \cdot 2 = 1$$

$$H(Y|X_3) = \dots = 1$$

$$\Rightarrow j = \arg \max_{i \in \{1, 2, 3\}} G(S, i) = 1 //$$

$$G(S, 1) = 1 - 0.69 = 0.31$$

the first split is by x_1 and we have two groups

$$S_0 = \{(001), 0\}$$

$$S_1 = \{(111), 1, (100)1, (111)0\}$$

for S_0 , it doesn't matter which $j \in \{2, 3\}$ we choose for the next split and the post split error will be 0 anyway.

we examine S_1 then

$$H(Y) = -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3}$$

$$H(Y|x_2) = H(Y|x_2=1)P(x_2=1) + H(Y|x_2=2)P(x_2=2)$$

$$\frac{2}{3} \left(-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right) + \frac{1}{3} \left(-1 \log 1 - 0 \log 0 \right) = \frac{2}{3}$$

$$H(Y|x_3) = \frac{2}{3} \left(-\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} \right) = \frac{2}{3}$$

$$\Rightarrow G(S_1, 2) = G(S_1, 3)$$

\Rightarrow it does not matter if we split by x_2 or x_3 . and we are going to mistake on 1 sample anyway

if we split by x_2 we have

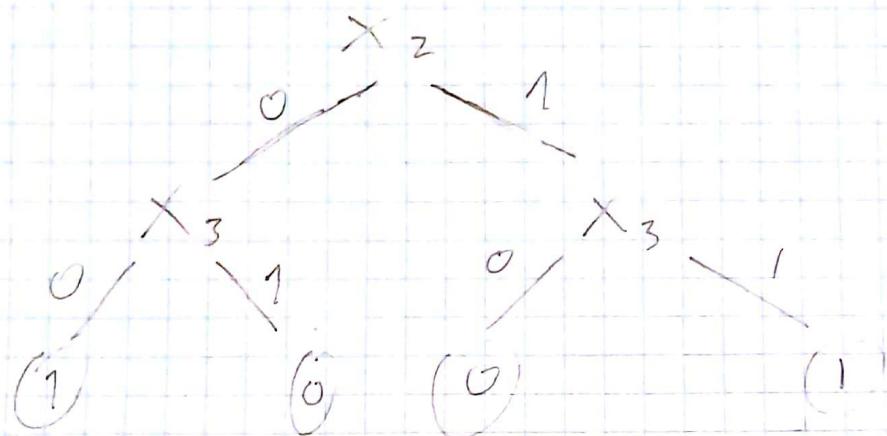
$$S_0 = \{(100), 1\} \rightarrow \text{no mistakes on this set}$$

$S_1 = \{(111), 1, (111), 0\} \rightarrow$ we have to choose a label by the majority now so we have 1 mistake guaranteed

similar for x_3 split $\Rightarrow 1$ errors out

$$\Rightarrow \text{error} = \frac{1}{4} \text{ of } n \text{ samples}$$

(5) we can see that a 0-error tree
of 2 splits can be



$$\textcircled{3} \quad (a) \quad \alpha_t = 0.5 \ln\left(\frac{1-\varepsilon_t}{\varepsilon_t}\right) = \ln\left(\frac{1-\varepsilon_t}{\varepsilon_t}\right)^{\frac{1}{2}}$$

$$\Rightarrow \varepsilon_t e^{\alpha_t} = \varepsilon_t (1-\varepsilon_t)^{\frac{1}{2}} = \sqrt{\varepsilon_t(1-\varepsilon_t)} \quad *$$

$$e^{\alpha_t} = \sqrt{(1-\varepsilon_t)}$$

$$\varepsilon_t e^{\alpha_t} = (1-\varepsilon_t) e^{-\alpha_t} \quad /$$

$$h_j = h(x_j)$$

$$\Rightarrow \sum_j D_t(x_j) e^{-\alpha_t y_j} h_j = \underbrace{\sum_{\text{mistakes}} D_t(x_j) e^{\alpha_t}}_{\varepsilon_t} + \underbrace{\sum_{\text{correct labeling}} D_t(x_j) e^{-\alpha_t}}_{1-\varepsilon_t}$$

$$= \underbrace{e^{\alpha_t} \sum_{\text{mistakes}} D_t(x_j)}_{\varepsilon_t} + \underbrace{e^{-\alpha_t} \sum_{\text{correct}} D_t(x_j)}_{1-\varepsilon_t}$$

$$e^{\alpha_t} \varepsilon_t + e^{-\alpha_t} (1-\varepsilon_t) = 2\sqrt{\varepsilon_t(1-\varepsilon_t)} \quad /$$

$$\textcircled{3} \quad (b) \quad \Pr_{x \sim D_{T+1}}(h_{T+1}(x) \neq y) = \sum_j D_{T+1}(x_j) I(h_j \neq y_j)$$

$$= \frac{\sum_j D_T e^{-\alpha_T y_j} h_j I(h_j \neq y_j)}{\sum_j D_T e^{-\alpha_T y_j} h_j} = \frac{\sum_j D_T e^{-\alpha_T y_j} h_j I(h_j \neq y_j)}{2\sqrt{\varepsilon_t(1-\varepsilon_t)}}$$

definition
 $\theta = D_{T+1}$

$$\frac{1}{2\sqrt{\varepsilon_t(1-\varepsilon_t)}} \sum_{j: h_j \neq y_j} D_T(x_j) e^{\alpha_t} = \frac{e^{\alpha_t}}{2\sqrt{\varepsilon_t(1-\varepsilon_t)}} \underbrace{\sum_{j=1}^n D_T(x_j) I(h_j \neq y_j)}_{\varepsilon_t} \quad /$$

$$= \frac{\varepsilon_t e^{\alpha_t}}{2\sqrt{\varepsilon_t(1-\varepsilon_t)}} = \frac{1}{2} \quad /$$

(c) Each h_i is a weak learner.

$$\epsilon_{D,S}(h_i) \leq \frac{1}{2} - \delta \quad 0 < \delta < \frac{1}{2}$$

but we showed that

$$\epsilon_{D_{T+1}}(h_{T+1}) = \frac{1}{2}$$

and

$$\epsilon_{D_{T+1},S}(h_{T+1}) \leq \frac{1}{2} - \delta < \frac{1}{2} = \epsilon_{D_{T+1}}(h_T)$$

$$\Rightarrow h_{T+1} \neq h_T$$

(d)

$$Z_T = \sum_i D_T(x_i) e^{-\alpha_T h_i(y_i)}$$

$$\frac{\partial Z_T}{\partial \alpha_T} = - \sum_i h_i(y_i) D_T(x_i) e^{-\alpha_T h_i(y_i)}$$

$$= \sum_{\text{mistakes}} D_T(x_i) e^{\alpha_T} - \sum_{\text{correct}} D_T(x_i) e^{-\alpha_T}$$

$$= e^{\alpha_T} \epsilon_T - e^{-\alpha_T} (1 - \epsilon_T) \stackrel{!}{=} 0$$

$$e^{2\alpha_T} = \frac{1 - \epsilon_T}{\epsilon_T} \Rightarrow \alpha_T = \frac{1}{2} \ln \left| \frac{1 - \epsilon_T}{\epsilon_T} \right|$$

by deriving again we can convince ourselves
that it is a minimum.

$$\textcircled{1} \quad S = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$y_i \cdot \sum_{j=1}^k \alpha_j h_j(x_i) \geq \gamma \quad \forall i$$

$$\gamma > 0, \alpha_j \geq 0 \quad \forall j, \sum_{j=1}^k \alpha_j = 1$$

$$\textcircled{2} \quad E(y_i \sum_{j=1}^k \alpha_j h_j(x_i)) \geq E(\gamma) = \gamma$$

$$\underset{i \in D}{E}\left(y_i \sum_j \alpha_j h_j(x_i)\right) = \sum_{i=1}^n D_i y_i \sum_{j=1}^k \alpha_j h_j(x_i)$$

$$= \sum_{j=1}^k \underbrace{\sum_{i=1}^n D_i y_i}_{\equiv g_j} h_j(x_i) = \sum_{j=1}^k \alpha_j g_j$$

$$= E(g_j) \geq \gamma$$

$j \text{ over } \{1, \dots, k\}$

\Rightarrow There exist $j^* \in \{0, \dots, k\}$ for which

$$g_{j^*} \geq \gamma$$

$$g_{j^*} = \sum_{i=1}^n D_i y_i h_{j^*}(x_i) = \sum_{i: y_i \neq h_{j^*}(x_i)} D_i y_i h_{j^*}(x_i) + \sum_{i: y_i = h_{j^*}(x_i)} D_i y_i h_{j^*}(x_i)$$

$$= - \sum_{\text{mistakes}} D_i + \sum_{\text{correct}} D_i = -P_{\text{mis}}(h_{j^*}(x_i) \neq y_i) - (1 - P(h_{j^*}(x_i) \neq y_i))$$

$$\Rightarrow -2P_{\text{ind}}(h_j(x_i) \neq y_i) + 1 \geq \gamma$$

$$\Rightarrow P_{\text{ind}}(h_j(x_i) \neq y_i) \leq \frac{1}{2} - \frac{\gamma}{2}$$

~~$j \in [0, k]$~~

(b)

$$S = \{(x_1, y_1), \dots, (x_n, y_n)\} \subseteq \mathbb{R}^d \times \{-1, 1\}$$

S realizable by a d-dim hyper rectangle classifier:

there exists a d-dim rectangle s.t.

$$y_i = 1 \text{ iff } x_i \in [a_1, b_1] \times \dots \times [a_n, b_n]$$

~~the rectangle~~

H is the class

$$h_{+\theta}(x) = \begin{cases} 1 & x(j) \leq \theta \\ -1 & x(j) > \theta \end{cases} \quad h_{-\theta}(x) = \begin{cases} 1 & x(j) \geq \theta \\ -1 & x(j) < \theta \end{cases}$$

coordinate j of some vector x

we set $|H| = nd - 1 \Rightarrow |H|$ hypotheses:

$$h_1 = h_{+\alpha_1}, h_2 = h_{-\beta_1}, \dots \text{ and so on}$$

those capture

The 1st coordinate of
 $x \rightarrow x_i(1) \in [a_1, b_1]$

for all dimensions

\Rightarrow overall $2d$ hypotheses.

and $2d - 1$ more of the form $h_{+\infty}$



now we set α_j (α 's in section a) $\forall j$:

$$\alpha_j = \frac{1}{nd-1} \quad \gamma = \frac{1}{nd-1}$$

(uniform distribution over hypothesis)

$$y \sum_{j=1}^k \alpha_j h_j(x(j)) \geq \gamma$$

inside rectangle:

Assume that for all $m = 1 \dots n$ $y_m = 1$

$$\Rightarrow x_m \in [a_1, b_1] \times \dots \times [a_d, b_d]$$

$$\Rightarrow j = 1 \dots d \quad h_j(x_m) = 1$$

$$j = d+1 \dots nd-1 \quad h_j(x_m) = -1$$

$$\Rightarrow y_m \sum_{j=1}^{nd-1} \alpha_j h_j(x_m) =$$

$$= \sum_{j=1}^{d-1} \frac{1}{nd-1} - \sum_{j=d+1}^{nd-1} \frac{1}{nd-1} = \frac{2d - 2d-1}{nd-1} = \frac{1}{nd-1} = \gamma$$

outside rectangle:

Assume $y_m = 1 \dots n$ $y_m(i) = -1$

$$x_m \notin [a_1, b_1] \times \dots \times [a_d, b_d]$$

\Rightarrow There exists p s.t. $x_m(p) \notin [a_p, b_p]$
($p \in \{1, d\}$)

$$\Rightarrow y_m \sum_{j=1}^{nd-1} \alpha_j h_j(x_m) = - \sum_{j \in \{d+1\}/p} \frac{1}{nd-1} h_j(x_m) - \underbrace{\frac{1}{nd-1} h_p(x_m)}_{=-1} - \sum_{j=d+1}^{nd-1} \frac{1}{nd-1} h_j(x_m)$$

$$\Rightarrow h_{d-1}$$

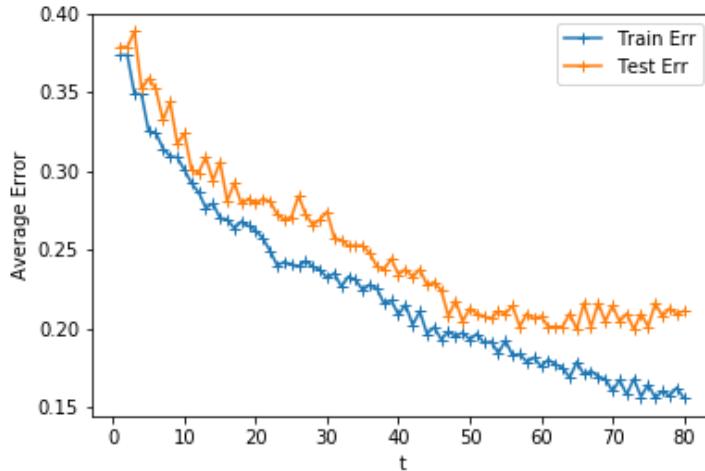
$$y_m \sum_{j=1}^{2d} \alpha_j h_j(x_m) = - \sum_{\substack{j \in [2d] \setminus p}} \frac{1}{nd-1} h_j(x_m) + \frac{1}{nd-1} + \frac{2d-1}{nd-1}$$

$$= - \sum_{\substack{j \in [2d] \setminus p}} \frac{1}{nd-1} h_j(x_m) + \frac{2d}{nd-1} \geq - \sum_{\substack{j \in [2d] \setminus p}} \frac{1}{nd-1} + \frac{2d}{nd-1} = \frac{1}{nd-1} = \gamma$$

$$\bullet \Rightarrow y_m \sum_{j=1}^{16} \alpha_j h_j(x_m) = \gamma \geq \gamma \quad \cancel{\text{✓}}$$

Programming assignment

- a. Plot of error (of classifier defined as a linear combination of weak learner hypotheses) on train data and test data with respect to t:



Both errors drop with the number of hypotheses that compose the classifier. It seems as if the test error has reached the optimum that can be achieved with the weak learners from our class.

- b. List of the first 10 words and their classification rule:

h0: word=**bad**: classify as 1 if word occur ≤ 0.5
h1: word=**life**: classify as -1 if word occur ≤ 0.5
h2: word=**many**: classify as -1 if word occur ≤ 0.5
h3: word=**worst**: classify as 1 if word occur ≤ 0.5
h4: word=**performances**: classify as -1 if word occur ≤ 0.5
h5: word=**plot**: classify as 1 if word occur ≤ 0.5
h6: word=**great**: classify as -1 if word occur ≤ 0.5
h7: word=**nothing**: classify as 1 if word occur ≤ 0.5
h8: word=**fun**: classify as -1 if word occur ≤ 0.5
h9: word=**script**: classify as 1 if word occur ≤ 0.5

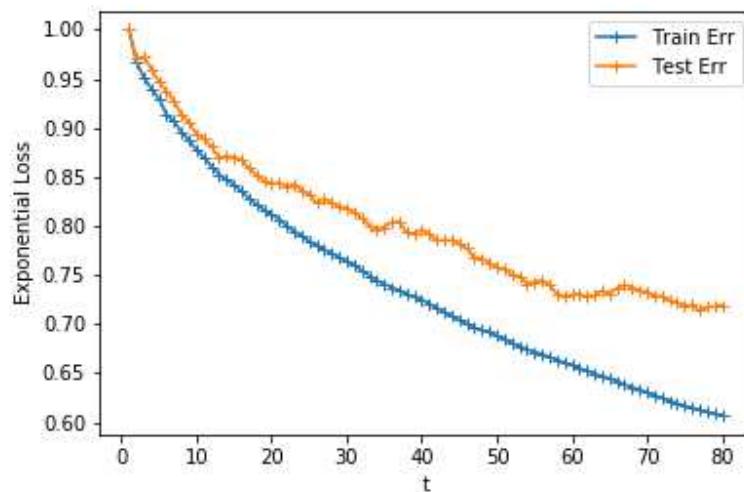
It seems that we have few words that are expected to be a good generalization for classifying reviews such as “bad”, “worst”, “great”. Intuitively, the first two would appear in bad reviews and the third in good reviews – note that this is indeed the rules we have for the hypotheses represented by the words.

We also have other words that are not expected to be related to a specific type of reviews, like “life”, “many”, “plot”. One explanation for the presence of such hypotheses is that after we get some hypothesis h' by our ERM, the next step in adaboost will focus on the samples that h' classified incorrectly. The next hypothesis will likely to regard a word that appears in those misclassified samples, just so it could classify correctly as many of them possible. Those words, sometimes are not good representatives of the whole data and are not intuitively good heuristics for reviews in general.

As an example, consider the review - “this movie is disappointing, unlike many other movies of this director that I thought was good”. The word “good”, if chosen as a classifier for good

reviews if they contain instances of it, will obviously mistake on that review. The next stage may choose the word “many” as a classifier for bad reviews just so this sample will then be classified correctly.

c. Average exponential loss of with respect to t



Loss of g on train data and test data drops with t (in other words – drops with the number of hypotheses achieved by adaboost). This is not surprising – we saw in recitation 10 that running adaboost is equivalent to minimizing the average exponential loss (up to a constant).

Also, we note that as in section a it seemed as if after enough iterations we achieved optimal error on the test set, here as well, the loss reach saturation after a while.