

Weak vs. Strong Learning

Def. Algorithm A weakly "PAC learns" concept class \mathcal{C}

if $\forall c \in \mathcal{C}$ + \forall dists \mathcal{D} $\exists \delta > 0$

$\forall \epsilon, \delta > 0$ ($\delta = \frac{1}{4}$ or $\frac{1}{n^2}$ doesn't affect)

with prob $\geq 1 - \delta$
given examples of c

A outputs h s.t. $\Pr_{\mathcal{D}} [h(x) \neq c(x)] \leq \frac{1}{2} - \frac{\delta}{2}$
 \uparrow
 advantage

It was conjectured that distribution free weak learning
was really weaker but surprise!

can "boost" a weak learner

Thm if \mathcal{C} can be weakly learned on
any dist \mathcal{D} then \mathcal{C} can be
(strongly) learned. \Leftarrow i.e. $\forall \epsilon, \beta_0$

Will prove for case of $\mathcal{D}_0 = \mathcal{U}$

Applications

1) "Theoretical"

- Unif dist Algorithms for poly term DNF
weight w - poly threshold fctns

} low degree
alg doesn't
work well

(Boosting + KM)

- Ave case vs. worst case complexity

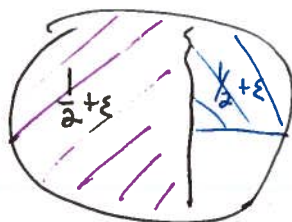
2) practical - Boosting
Freund-SchapireGood & Bad Ideas

- 1) simulate weak learner several times on
same distribution & take majority answer
-or-
best answer

gives better confidence

but doesn't reduce error, what if always get same answer?

- 2) filter out examples on which current hypothesis
does well & run weak learner on part where you
do badly.

Problem: given a new
example, how do you
know which section it
is in?

- 3) **Keep** some samples on which you are ok in filtering
 always use **majority vote** on all previous hypotheses
 to predict value of new samples

history: Schapire, Freund-Schapire, Impagliazzo-Servedio, Klivans

Filtering Procedures

- decide which samples to keep, which to throw out
- samples on which so far you guess correctly \leftarrow need for check future hypot
- samples on which so far you guess incorrectly \leftarrow need to improve on these

The setting

- Given labelled examples
 $(x_1, f(x_1)), (x_2, f(x_2)), \dots$

$$x_i \in \mathcal{X}$$

$$f \in \mathcal{C}$$

- Given weak learning alg WL which weakly learns (advantage $\frac{\epsilon}{2}$) on any dist \mathcal{D}

Boosting Algorithm

• Stage 0 (Initialize)

$$D_0 \leftarrow D$$

run WL on D_0 to generate (whp)

$$C_1 \text{ s.t. } \Pr_{D_0} [f(x) = C_1(x)] \geq \frac{1}{2} + \gamma/2$$

• For $i = 1 \dots T = O(\frac{1}{\gamma^2 \epsilon})$ stages, stage i : (can stop if Majority($C_1 \dots C_i$) correct on $\geq 1-\epsilon$ inp)

(1) Construct D_i via "filtering procedure":

{ favor pts on which maj of $C_1 \dots C_i$ don't do well
but also keep some other points }

Will specify soon

(2) run WL on examples from D_i to output

$$C_{i+1} \text{ s.t. } \Pr_{D_i} [f(x) = C_{i+1}(x)] \geq \frac{1}{2} + \frac{\gamma}{2}$$

• output $C = \text{MAJ}(C_1 \dots C_T)$

Filtering procedure

Given new example $x, f(x)$ from example oracle

• if majority of $C_1 \dots C_i$ wrong, keep it
ie. $\geq \frac{i}{2}$

• if large majority right, then discard

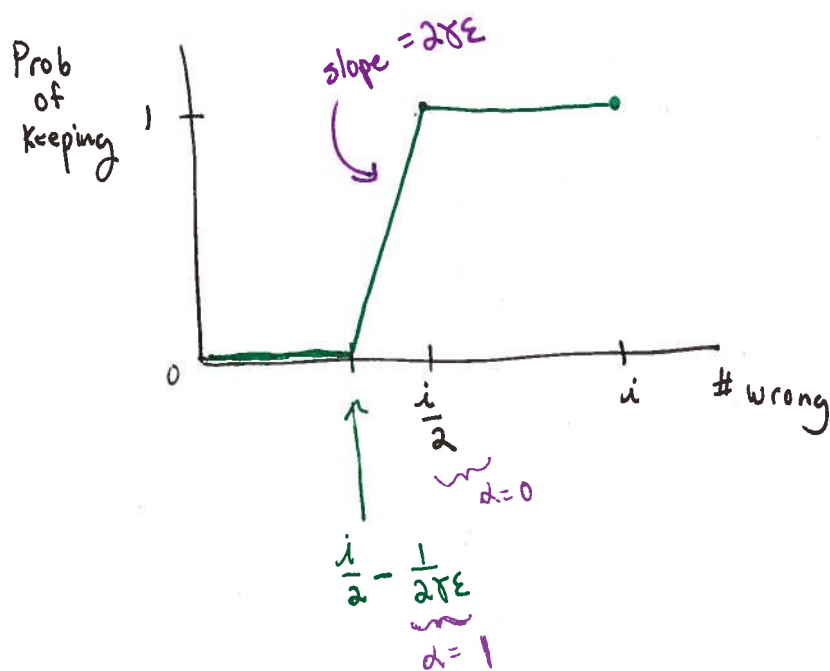
$$\text{ie. } \# \text{ right} - \# \text{ wrong} > \frac{1}{\gamma \epsilon}$$

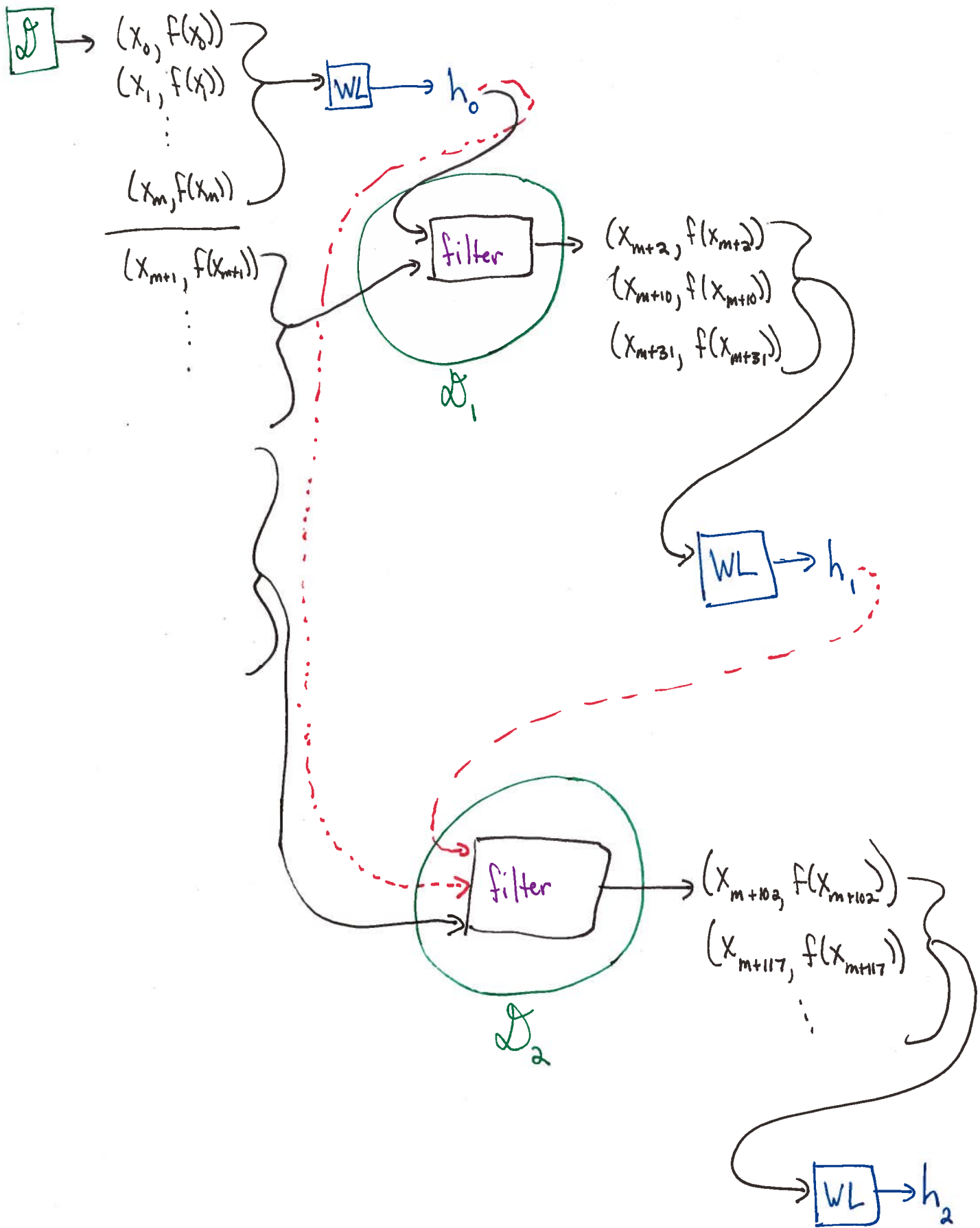
$$\text{or } \# \text{ wrong} \leq \frac{i}{2} - \frac{1}{2\gamma \epsilon}$$

• else $\# \text{ right} - \# \text{ wrong} = \frac{\alpha}{\gamma \epsilon}$ for $0 < \alpha < 1$

$$\# \text{ wrong} - \# \text{ right} = \frac{-\alpha}{\gamma \epsilon}$$

So keep with prob = $1 - \alpha$





Need to show:

1) Output is has nontrivial agreement with f

2) # samples needed not too bad

why could it be bad?
 if throw out lots of samples, might
 need to wait a long time before WL
 can give an output, too many samples then
 but if throw out too many samples then
 you already have a good hypothesis!

will stop if $\text{Maj}(C_1, \dots, C_i)$ correct on $\geq 1-\epsilon$ fraction of inputs

ow. $\text{Maj}(C_1, \dots, C_i)$ incorrect on $> \epsilon$ fraction

so filtering procedure outputs sample with prob $\geq \epsilon$ (+ in expectation, every $1/\epsilon$ samples of \mathcal{D} at least one makes it thru the filtering system)

\Rightarrow filtering slows down sample collection by $\leq O(1/\epsilon)$

Then filter outputs new sample in time $\leq \frac{1}{\text{error}[\text{Maj}(C_1, \dots, C_i)]}$

So lets focus on ①

Notation

$$R_c(x) = \begin{cases} +1 & \text{if } f(x) = c(x) \\ -1 & \text{if } f(x) \neq c(x) \end{cases}$$

"is c correct on x?"

$$N_i(x) = \sum_{1 \leq j \leq i} R_{c_j}(x)$$

after iteration i_j
how many c's correct?
(#right - #wrong)

$$M_i(x) = \begin{cases} 1 & \text{if } N_i(x) \leq 0 \\ 0 & \text{if } N_i(x) \geq \frac{1}{\epsilon} \\ 1 - \epsilon \cdot N_i(x) & \text{o.w.} \end{cases}$$

prob of keeping x
in filtering
(after stage i)
note - all "wrong" x included
in M
also some "right" x included

Note that new distribution on samples is proportional to M_i :

$$D_{M_i}(x) = \frac{M_i(x) D_0(x)}{\sum_y M_i(y) D_0(y)}$$

since D_0 is assumed to be uniform, we can drop this & note $D_{M_i}(x) = \frac{M_i(x)}{\sum_y M_i(y)}$

$\sum_y M_i(y)$ includes all "wrong" y but also y for which maj that isn't are correct } upper bounds # wrong y

How correct are we w.r.t. D_{M_i} ?

$$\text{Adv}_c(M_i) = \sum_x R_c(x) M_i(x)$$

"Advantage" of c on M
 $\sim \Pr[\text{correct}] - \Pr[\text{incorrect}]$
 $= 2 \cdot \Pr[\text{correct}] - 1$

$$\Pr_{x \in D_{M_i}} [c(x) = f(x)] = \frac{1}{2} + \frac{\text{Adv}_c(M_i)}{2 \cdot \sum_x M_i(x)}$$

Note:

$$\text{if } \sum M_i(x) \geq \epsilon 2^n$$

$$\text{Adv}_c(M_i) \geq \gamma \cdot \epsilon \cdot 2^n$$

convert claim about WL \Rightarrow claim about advantage
 i.e. if have γ advantage on output of i
 + still almost wrong on lots of inputs
 then new advantage is pretty good
 if not, then you are done

Begin Proof

For input x

$$\text{let } A_i(x) \leftarrow \sum_{0 \leq j \leq i-1} R_{c_{j+1}}(x) M_j(x)$$

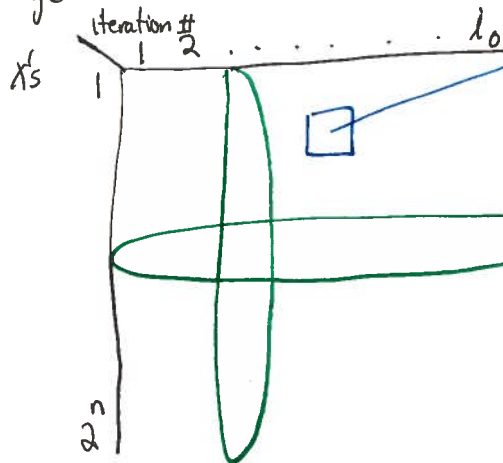
strange -
 indices don't ma
 c_1, \dots, c_j define D_j
 but c_{j+1} learned A
 WL on D_j

Claim $A_i(x) \leq \frac{1}{\epsilon \gamma} + \frac{\epsilon \gamma}{2} \cdot i$

- bounds advantage per input
- only helps after $\frac{1}{\epsilon \gamma}$ rounds

Plan for use of claim:

Consider large matrix:



weighted version of whether c_{j+1} correct on x

(k,j) th Entry: $R_{c_{j+1}}(x_k) M_j(x_k)$

$$x\text{'s row sum} = \sum_{0 \leq j \leq i_0} R_{c_{j+1}}(x) M_j(x) = A_{i_0+1}(x)$$

$$j\text{th col sum} = \sum_x R_{c_{j+1}}(x) M_j(x)$$

$$= \text{Adv}_{c_{j+1}}(M_j) \leftarrow \begin{cases} \geq \frac{\epsilon \gamma}{2} & \text{if } \text{else algorithm } s \end{cases}$$

Goal: lower/upper bound average entry in matrix

lower bound:

lower bound average entry in column via

- correctness of WL

- fact that algorithm proceeds

\Rightarrow lots of error

$\Rightarrow \sum_x M_j(x)$ big

\Rightarrow lots of progress in WL
in absolute terms

upper bound:

upper bound rows via claim

- if advantage is too much, lose measure

this is where majority rule
& weighting scheme is used

More details:

Assume claim, prove theorem:

Assume haven't terminated at $i_0 + 1$ th stage

- so error $(C_{i_0}) \geq \epsilon$

$$\sum_x M_{i_0}(x) \geq \epsilon 2^n$$

Claim \Rightarrow

$$\sum_x A_{i_0+1}(x) = \sum_x \sum_{0 \leq j \leq i_0} R_{C_{j+1}}(x) M_j(x) \quad \text{def of } A_{i_0+1}$$

$$= \sum_{0 \leq j \leq i_0} \text{Adv}_{C_{j+1}}(M_j) \quad \text{def of } \text{Adv}_{C_{j+1}}$$

$$\geq (\gamma \epsilon 2^n) (i_0 + 1)$$

From "note"

$$+ \sum_x A_{i_0+1}(x) \leq \sum_x \left(\frac{1}{\epsilon \gamma} + \frac{\epsilon \gamma}{2} \cdot (i_0 + 1) \right) \quad \text{claim}$$

$$= 2^n \left(\frac{1}{\epsilon \gamma} + \frac{\epsilon \gamma}{2} (i_0 + 1) \right)$$

putting together:

$$(\epsilon \gamma) (i_0 + 1) \leq \frac{1}{\epsilon \gamma} + \frac{\epsilon \gamma}{2} (i_0 + 1)$$

$$\text{so } \frac{\epsilon \gamma}{2} (i_0 + 1) \leq \frac{1}{\epsilon \gamma} \Rightarrow i_0 \leq \frac{2}{\epsilon^2 \gamma^2} - 1$$

Proof of claim:

Question: how can an input x add to cumulative advantage throughout algorithm?

Observations:

- if algorithm's hypotheses $c_1 \dots c_i$ are overwhelmingly correct on x , then not at all because x gets measure 0
- if algorithm's hypotheses are doing badly (mostly wrong) then not at all because they decrease advantage

• Main Issue:

can wander in middle -
majority correct but not large majority
increase advantage so have positive measure

need to bound this case.

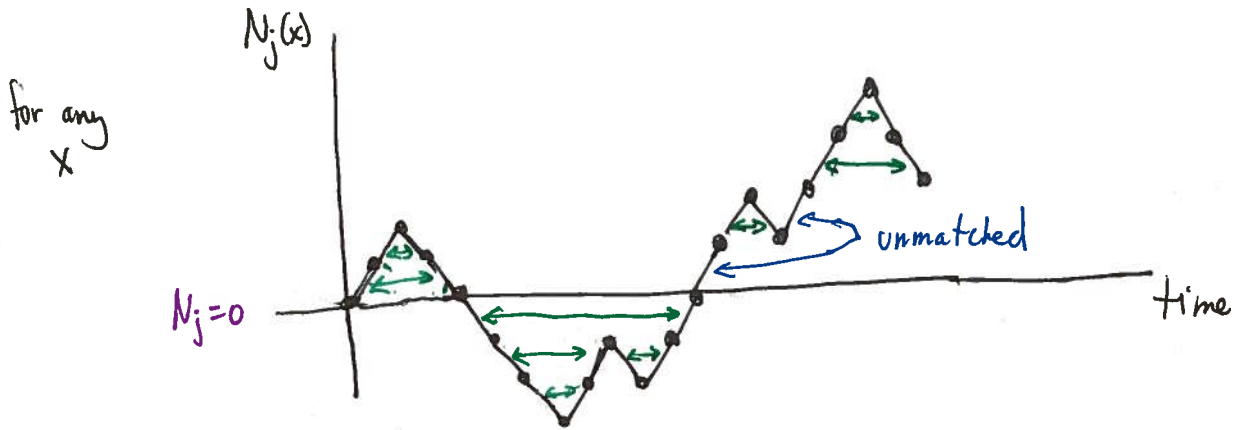
Proof of Claim

First, strange but obvious fact:

Fact "elevator argument"

If one spends any amount of time in an elevator, then no matter what sequence of buttons pushed, one ascends from k^{th} to $k+1^{\text{st}}$ floor at most one more time than one descends from the $k+1^{\text{st}}$ to k^{th} floor.
 (analogous for negative floors $-k$ & $-(k+1)$)

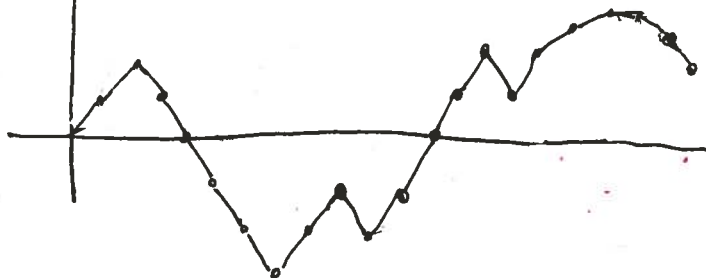
Proof by picture:



match transitions going up with those going down on same level (as in parentheses)

but what is behavior of $\sum_{j \leq k} R_{c_{j+1}}(x) M_j(x)$?

$$\sum_{j \leq k} R_{c_{j+1}}(x) M_j(x)$$



± 1
 $\in [0, 1]$
 $\Rightarrow |\text{slope}| \leq 1$ (in fact, $\leq 2\epsilon$)
 + same sign as $N_j(x)$

Recall: $A_i \equiv \sum_{0 \leq j \leq i-1} R_{c_{j+1}}(x) M_j(x)$

Matching:

For $k \geq 0$:

match $a = j$ s.t. $N_j(x) = k$ + $N_{j+1}(x) = k+1$

with $b = j'$ s.t. $N_{j'+1}(x) = k+1$ + $N_{j'+2}(x) = k$

For $k < 0$: analogously match $-k$ to $-(k+1)$
 with $-(k+1)$ to $-k$

use N_i 's to create matching. the corresponding R_{c_i} 's will be such that one is $+1$ and one is -1

For each matched pair:

Will bound contribution from matched pairs
 by $\epsilon \delta$ per pair using bound on slope
 (and total of $\frac{\epsilon \delta i}{2}$)

(for each matched pair (a, b) cont.)

just by assumption that $R_{c_{a+1}}(x) = +1$ + $R_{c_{b+1}}(x) = -$

$$\underbrace{R_{c_{a+1}}(x)}_{\substack{+1 \\ \text{elevator} \\ \text{going up}}} \underbrace{M_a(x)}_{N_a(x)=k} + \underbrace{R_{c_{b+1}}(x)}_{\substack{- \\ \text{elevator} \\ \text{going down}}} \underbrace{M_b(x)}_{N_b(x)=k+1} = M_a(x) - M_b(x)$$

if $0 \leq k \leq \frac{1}{\epsilon\gamma}$ or $0 \leq k+1 \leq \frac{1}{\epsilon\gamma}$

$$\begin{aligned}
 \text{then } & \underbrace{M_a(x)} - \underbrace{M_b(x)} \\
 &= (1 - \epsilon\gamma N_a(x)) - (1 - \epsilon\gamma N_b(x)) \\
 &= \cancel{1} - \cancel{\epsilon\gamma k} - \cancel{1} + \epsilon\gamma(k+1) \\
 &= \epsilon\gamma
 \end{aligned}$$

Contribution of matched edges in the "low slope" area

$$\text{else } M_a(x) - M_b(x) = \begin{cases} 1-1 \\ \text{or} \\ 0-0 \end{cases} = 0$$

Contribution of matched edges in the "slope=1" area

\therefore each pair $\leq \frac{i}{2}$ pairs contributes $\leq \epsilon\gamma$ to sum $\left. \begin{matrix} \} \\ \} \end{matrix} \right\} \leq \frac{i}{2} \cdot \epsilon\gamma$ total contribution

Contribution from unmatched edges:

either all unmatched N_i 's have negative steps
 or all have positive steps

if all negative:

R_{c_j} 's all -1

M_j 's all $\in [0, 1]$

\therefore contribution of $R_{G_H}(x) M_j(x) < 0$

if all positive:

if unmatched N_i 's in $[0, \frac{1}{\epsilon\gamma}]$

- for each $M_j \in [0, 1]$, contribution of

$$R_{c_{j+1}} M_j(x) \leq 1$$

- at most $\frac{1}{\epsilon\gamma}$ of these

$$\Rightarrow \text{total contribution} \leq \frac{1}{\epsilon\gamma}$$

if unmatched N_i in $[\frac{1}{\epsilon\gamma}, \infty]$

then $M_j = 0$

\Rightarrow total contribution = 0

$$\therefore \text{Grand total} \leq \frac{1}{2} \cdot \gamma \epsilon \cdot i + \frac{1}{\epsilon\gamma}$$

