

# Contents

<b>1</b>	<b>Big picture</b>	<b>1</b>
1.1	Boosting . . . . .	2
<b>2</b>		<b>6</b>
<b>3</b>	<b>Improved boosting algorithms and implications</b>	<b>12</b>
<b>4</b>	<b>2017-3-2</b>	<b>16</b>
	Wiki page: <a href="https://simons-institute.github.io/pseudorandomness/groups/modelstructure.html">https://simons-institute.github.io/pseudorandomness/groups/modelstructure.html</a>	

## 1 Big picture

I'll talk about things which have different names in different fields. You probably know them, but don't know they are the same or related ideas that come up in different fields.

	Boosting	Hard-core lemma	Dense model theorem	Weak regularity	?
Area	ML	CC, Derandomization	Additive combinatorics, CC	Graph theory	
Credit	Shapiro, Freund-Schapire	Impagliazzo, Holenstein	Green-Tao, Barak-Shaltiel-Wigderson	Szemerédi, Frieze-Kannan	
Get	Circuit computing $f$ $1 - \delta$ of the time	"	Proof that set isn't $\delta$ -dense	"	
Unless	Weak learner fails on distribution of density $\Omega(\delta)$	Hard-core distribution	$\Omega(\delta)$ -dense "model" indistinguishable from set	A model succinctly describing set	
Algorithm needed	Weak learner	"	Distinguisher	"	

Implications are due to...

1. Boosting  $\implies$  Hard-core: Klivans and Servedio.
2. Hard-core  $\implies$  Dense model: Impagliazzo
3. Dense model  $\implies$  Weak regularity: Trevisan-Tulsiani-Vadhan, Reingold-Trevisan-Tulsiani-Vadhan

4. Weak regularity  $\implies$  boosting: Trevisan-Tulsiani-Vadhan

If set looks like might be dense, either succinct reason (proof) not  $\delta$ -dense, or dense model (reason why it looks big to certain class of sets)

Hard-core lemma implies dense model theorem.

FK simpler argument with better quantitative bounds, weaker conclusion. Generalized to sparse graphs without giveaway cut (cut prove not dense). Either you get a proof that the cut isn't dense, or get a model that succinctly describes all cuts.

Take a bunch of theorems we know are true and showing implications between them. Hope a lot to gain by looking at connections between them:

1. Versatility. Retrofit. Ex. there are many boosting algorithms. When you follow this progression, you get different quantitative and qualitative versions of dense model theorem and regularity.

Boosting is a particular proof for the min-max theorems.

2. Algorithmic and constructive:

There are nonconstructive versions using the min-max theorem for boosting, hard-core lemma, dense model theorem.

I care about algorithmic versions.

In ML we care about getting a function that computes a function much of the time.

On the other side, we're really after the distribution where the weak learner fails, the model that succinctly describes the set. We pay attention to do the reductions in an algorithmic, not just an existential way.

3. Dense model theorem  $\rightarrow$  learning. Take boosting technique and use it in an unsupervised way.
4. Generality: some things seem to be specific to a setting (density of graphs). Actually weak regularity doesn't have anything to do with graphs being dense. We can relativize it to subgraphs of any graph. You can look at subgraphs of expanders, bipartite graphs, etc. and plug it in the same machinery. If you want to look at spectral norms rather than cuts, you can plug into the same machinery.

Almost all arrows are trivial! Simple lines of high-school algebra.

Weak regularity does imply dense model theorem and hard-core lemma.

## 1.1 Boosting

First we discuss the PAC learning model.

Let  $U$  be the universe—big set of possible inputs. There is a boolean function  $f$  and we get points  $(x, f(x))$  where  $x \sim U$ . The algorithm can request any number of such points. Our goal is to find a hypothesis  $h$  such that

$$\mathbb{P}_{x \sim U}[h(x) = f(x)] \geq 1 - \delta.$$

How do we get a hypothesis  $h$ ? We assume we have a weak learner, which receives the same kind of input, but perhaps from some other distribution  $\mu$ , and gets something correlated with the right answer.

$$\mathbb{P}_{x \sim \mu}[h_i(x) = f(x)] \geq \frac{1}{2} + \varepsilon.$$

(It has to work on any distribution  $\mu$  satisfying some assumptions.)

How do we use the weak learner? Give a routine that subsamples the distribution  $U$  to pass to  $\mu$ . We let  $\mu : U \rightarrow [0, 1]$ . Keep something in  $U$  and keeping it with probability in  $[0, 1]$ , else throw it back (rejection sampling). This induces a probability distribution

$$D_\mu(x) = \frac{\mu(x)}{\sum_{x' \in U} \mu(x')}.$$

Let  $|\mu| = \mathbb{E}_{x \in U} \mu(x)$ . Use this as a notion of density. Condition:

1. We will only run the weak learner on distributions of about the same density,  $|\mu| \geq \Omega(\delta)$ .

We don't want to run it on something too puny, because the time to simulate the distribution is inversely proportional to the density.

2. Each  $h_i \in T$ . Ex.  $T$  is linear separators.

$$F_k T = \{f(h_1, \dots, h_k) : f \in F, h_1, \dots, h_k \in T\}.$$

It turns out that  $\mu \in F_k T$ : it's basically describable in the same class of functions.  $h \in F_k T$ .

Pick one large sample, run on it as much as you'd like. Then the issues of what is  $\mu$  disappears.

If you just care about ML we don't care about this details, but to go further, we do.

The first boosting algorithm I'll give is totally ridiculous from the ML point of view. For people who work on weak regularity on graphs this is the natural version and leads to the standard versions of results.

Suppose you just look at things information-theoretically. If I learn these hypotheses are relevant, how should I use them to predict the function? What's the best function that composes  $k$  of these things and gives a predictor for  $f$ ?

$h_1$  is either T or F, etc. This divides into  $2^n$  buckets. Once you are in a bucket what should you do? When given the next  $x$ , all you know is which bucket it's in.

To optimize, you predict the majority. Let  $p_B$  be the probability of the majority; then  $1 - p_B$  is the minority. Let  $H_i(x)$  be the majority answer in  $B$ .  $B$  is distributed according to  $B(x)$ . It takes  $2^i$  time to give a good estimate for the majority answer.

Let  $\delta_i = \mathbb{E}_{B \in B[x], x \in U} [1 - p_B]$ . If the majority  $\rightarrow 1$  in every bucket, then our error probability is small. While it's big, we still need to learn.

On the same biased distribution, get the same  $h_1$ . We want to make sure not boolean combination of previous  $h_i$ 's is a good predictor for the next distribution. If  $f(x) = \text{Maj}_B$ , then  $\mu(x) = \frac{1-p_B}{p_B}$  for  $x \in B$ .

What is  $|\mu|$ ?

$$\mathbb{E}_x |\mu| = \mathbb{E}_x \mu(x) = \mathbb{E}_B [\mathbb{E}_x [\mu(x) | x \in B]] \quad (1)$$

$$= \mathbb{E}[1 - p_B + 1 - p_B] = 2 \mathbb{E}_B [1 - p_B] = 2\delta_i. \quad (2)$$

If measure small, then small failure probability. If haven't succeeded, then  $\mu$  is large enough so we can continue learning on  $\mu$ .

Given previous functions, define candidate function. Either it has good success, or distribution has good density, and we continue. I need to show this process terminates within a reasonable amount of time.

We'll look at a potential function that measures how close we are to success. Our goal is to drive all  $p_B$ 's to 1.

Our potential function is

Flip a coin with expected probability.

$$\varphi_i = \mathbb{E}_{p_\varphi} p_1(1 - p_1) = \mathbb{E}_B p_1 - p_1^2 \quad (3)$$

$$= \mathbb{E}_B p_{B,1}(1 - p_{B,1}) \quad (4)$$

Land in  $B_1$  with probability  $q$ ,  $B_0$  with probability  $1 - q$ .  $f = 1$ ,  $p_{B,1} + \alpha_+$ ,  $p_{B,1} - \alpha_-$ .  
 $q\alpha_t = (1 - q)\alpha_-$ .

???

$$q(p_{B_1} + \alpha_+)^2 + (1 - q)(p_{B_1} - \alpha_-)^2 = (q + 1 - q)p_{B,1}^2 + (\alpha_t q - (1 - q)\alpha_- - 1)p_{B,1} + q\alpha_t^2 + (1 - q)\alpha_-^2. \quad (5)$$

$$\mathbb{E}_x [\mu(x)(-1)^{f(x) \oplus h_{i+1}(x)}] \geq \varepsilon |\mu| \quad (6)$$

$$= \mathbb{E}_{B \in B(x)} \mathbb{E}_{x \in B} \mu(x)(-1)^{f(x) \oplus h_{i+1}(x)} \quad (7)$$

$$q \left( \frac{1 - p_B}{p_B} (p_B + \alpha_t) - (1 - p_B + \alpha_-) \right) \quad (8)$$

$$(1 - q) \left( \frac{1 - p_B}{p_B} (1 - p_B + \alpha_-) - (p_B + \dots) \right) \quad (9)$$

Prob  $h$  on  $B$  is  $q$ . For  $f = 1$ ,  $q + \alpha_+$ . For  $f = -1$ ,  $q - \alpha_-$ . Then

$$p\alpha_+ = (1 - p)\alpha_-.$$

Then

$$\mathbb{E}_\mu \mu(x)(-1)^{f(x) \oplus h(x)} = \varepsilon_B = p \frac{1 - p}{p} (q + \alpha_+ - (1 - q - \alpha_+)) \dots \quad (10)$$

$$(1 - q - \alpha_-) - (1 - \alpha_-) \quad (11)$$

$$(1 - q)(2\alpha_+ + 2\alpha_-). \quad (12)$$

On the other hand the expectation

$$p_{B,1} = q \left( \frac{p(q + \alpha_+)}{q} \right)^2 - (1 - q) \left( \frac{p(1 - q - \alpha)}{1 - q} \right)^2 \quad (13)$$

$$= p^2 \left( q + \frac{2\alpha_+}{q} + \frac{\alpha_+^2}{q} \right) () \quad (14)$$

Try to make an energy increment argument.

It matters that the weak learner only looks at distributions with density  $\geq 2\delta$ . From boosting to hard-core:

**Lemma 1.1.** *Let  $f$  be any function on  $U$  and  $T$  any class of tests. Either there exists  $H \in F_k T$ ,  $H(x) = f(x)$  with probability  $1 - \delta$ , or there exists  $\mu$ ,  $|\mu| \geq 2\delta$ ,  $\forall h \in T$ ,*

$$\mathbb{P}_{x \in \mu}[h(x) = f(x)] \leq \frac{1}{2} + \varepsilon,$$

$$k = O\left(\frac{1}{\varepsilon^2 \delta^2}\right).$$

*Proof.* Let weak learner be exhaustive search. get sequence of distribution, all of  $\geq 2\delta$ . See if any of them have a good bias. If we always get something with bias  $\geq \frac{1}{2}$ , continue in the strong learner, until we get  $H(x) = f(x)$  with probability  $1 - \delta$ . If at some point our exhaustive search algorithm gets stuck, we have a distribution that's hard core.  $\square$

Region where very hard, indistinguishable from random.

Why have to keep track: we want size of hard-core to be what it should be.  $2\delta$  is what it should be.

A lot of times we want to compute one or the other. Then we need to use the boosting algorithm. By using the boosting algorithm we get something stronger. We get that the hardcore distribution doesn't have unlimited complexity. It's describable by the function we want to compute,  $\mu \in F_{k+1}(J \cup \{f\})$ .

Let  $S \subseteq U$ . We have some tests  $T$  that don't reveal that it's small. If big, then  $\delta$  fraction of whole thing.

**Definition 1.2:** Let  $S \subseteq U$ . Let  $\mathcal{T}$  be a class of tests. For  $T \in \mathcal{T}$ ,  $\mathbb{E}_{x \in S} T(x)$ .  $S$  is  $(\varepsilon, \delta)$ -pseudo-dense for  $\mathcal{T}$  if for all  $T \in \mathcal{T}$ ,  $T(U) \geq \delta T(S) - \varepsilon$ . This is written as  $S \sim_{\varepsilon, \delta} U$ .

One way to be pseudo-dense it to be dense. Another, one step removed, is that there's  $R$  that's indistinguishable from the whole distribution, and it's dense in  $R$ . and  $T(U)$  and  $T(R)$  are within  $\delta$  fraction of each other.

**Theorem 1.3** (Dense model theorem). *There exists  $\mu$ ,  $|\mu| \geq \delta$  and  $S \sim_{\varepsilon, \delta} D_\mu$ .*

*If  $S$  is  $(\varepsilon, \delta)$ -pseudodense vs  $F_k \mathcal{T}$ ,  $k = O\left(\frac{1}{\varepsilon^2 \delta^2}\right)$  then there exists  $\mu$ ,  $\mu \in F_k \mathcal{T}$ ,  $|\mu| \geq \delta - \varepsilon$ ,  $D_\mu \sim_{\mathcal{T}} S$  to within  $O(\varepsilon/\delta)$ .*

By going through boosting and hard-core, we get characterization of  $\mu$ , describable using original tests.

Hard function should be membership in  $S$ . This happens almost never, so we have to blow up membership to make it comparable to what it should be.

*Proof.* Let  $U'$  be:  $\delta'$  fraction is  $(0, x), x \in S$ , while  $1 - \delta'$  is  $(1, x), x \in U$ . Let  $\delta' = \frac{\delta}{1+\delta}$ . Let  $f(b, x) = b$ ,  $T((b, x)) = f(b, x)$  with probability  $1 - \delta'$ .

The tests are just applied to the  $x$  part.

$$\delta' T(S) + (1 - \delta')(1 - T(U)) = 1 - \delta' + \delta'(T(S)) - (1 - \delta')T(U) \quad (15)$$

$$= 1 - \delta' + \frac{1}{1 + \delta}(\delta T(S) - T(U)) \leq 1 - \delta' + \varepsilon. \quad (16)$$

This distribution has a hardcore subdistribution on which is as hard as a random coin flip. With  $\delta'$  probability, sample from  $S$ ,  $f = 1$ . with  $1 - \delta'$  probability, sample from  $U$ ,  $f = 0$ . The subset has size  $2\delta'$ , in order to be hard, half comes from  $S$  (all of it), half from  $U$ .  $f$  hard to predict means the 2 sides indistinguishable.

The model depends on  $k$  and  $f$ . Magically when we apply it we only want the part of the model where  $f = 0$ . It's definable with  $k$  things from  $\mathcal{T}$ . No membership oracle required.  $\square$

## 2

I'll start with motivation: one thing we could hope to get from this.

1. Let  $X$  be a set, e.g. a distribution of points in the square.

Let  $S$  be some distribution on points in  $X$ .

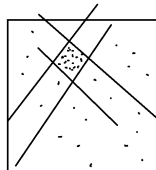


Let  $\mathcal{T}$  be a set of classifiers, ex. a set of half-planes.

Let  $\mathcal{F}_K \mathcal{T}$  be boolean functions on  $K$  functions in  $\mathcal{T}$ ; here, partitions into polygonal regions by  $k$  half-planes.

We want to pre-process the distribution to be able to answer queries in  $\mathcal{F}_K \mathcal{T}$ .

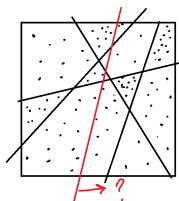
2. A violation of pseudo-density in this setting means there is a polygonal region with many more points from  $S$  than its volume, a “hot spot”.



$$\text{Area}(\text{region}) < \delta \mathbb{P}_S(\text{region}) - \varepsilon.$$

3. A model is a partition into polygonal regions, with a probability distribution on regions. A simple model is defined by at most  $k$  lines.

The property of a model is that we can estimate half-space probabilities (“given any half-space, what proportion of points of  $S$  are on one side of it?”) by treating the points as if uniform within regions.



4. The algorithmic requirement in order to process the points to answer queries in  $\mathcal{F}_K \mathcal{T}$  is: given a set of points sampled from  $S$ , and a set of points sampled from  $U$ , find a half-space that approximately maximizes the difference in probabilities for these two sets. The equivalent in boosting is a distinguishing algorithm.

Setting	Boosting	Hard-core measure	DMT/transference principle	Weak regularity
	WL: $ \mu_i  \geq 2\delta$ , $\mu_i = g(h_1, \dots, h_i, f)$ , $h_{i+1} \in \mathcal{T}$ , $k$ iterations	Hardcore measure: $\mu_k = g(h_1, \dots, h_k, f)$ , $ \mu_k  \geq 2\delta$	Model: $\mu_k = g(h_1, \dots, h_k, o)$ , $ \mu_k  \geq \delta$	
	SL: $H = G(h_1, \dots, h_k)$ , $\mathbb{P}[H = f] \geq 1 - \delta$	Violation of hardness: $H = G(h_1, \dots, h_k)$ , $\mathbb{P}[H = f] > 1 - \delta$	Violation of pseudo-density: $H = G(h_1, \dots, h_k)$ , $H(U) \leq \delta H(S) - \varepsilon$	
Assumption	WL never fails	Violation is impossible	Violation of pseudo-density is impossible	Actually dense
Conclusion	SL works	Hard-core measure exists, with same $k$ , $G, g$	Model exists	Model exists
Algorithmic	Weak learner requirement	Approximately optimal weak learner	Approximately optimal distinguisher	

Weak regularity: Same as above

Assumptions: Actually dense

Conclusion: model exists

Some comments:

1. Boosting: Note it's important that the  $\delta$  here is the same; many boosting algorithms meet this criterion.

The theorem says that “either weak learner fails or strong learner works.”

In boosting, we think of weak learner as never failing.

2. Hard-core measure lemma: The lemma says that either we can find hard-core measure, on which no function can compute the function  $f$  more than  $\frac{1}{2} + \delta$  of the time; or find a function computing  $f$   $1 - \delta$  of the time.

Here, we want to come up with the measure. Although the logical format is the same as boosting, here we assume that the violations never happen (there is no strong learner).

Every boosting algorithm gives hard-core measure lemma with the same parameters, and with exactly the same way of “gluing” the functions.

Sometime you care about computational complexity of  $G$  but not of  $g$ , or vice versa.

3. We can convert the hard-core measure theorem into the dense model theorem/transference principle (Tao and Ziegler).

Here, we have a distribution we're trying to model. Either the distribution has pseudo-density property—there isn't a violation that's definable from  $k$  different properties from hypothesis class, where violation means that the expected value is much smaller on  $U$  than on  $S$ —or we get a model of density  $\geq \delta$ .

Assuming that violation of pseudo-density does not happen, we get a model.

4. Weak regularity is just DMT except the distribution actually is dense. It's not so interesting that it has a dense model.

What we get is that the dense model you get is simple, definable in terms of a small number of basic hypotheses.

Sometimes we care about simplicity in the model, and sometimes simplicity in  $G$ .

5. Note the  $k$  is the same throughout. Reductions preserves  $k$ , and the functions  $h_i, G$ .

We don't only have the fact that boosting implies hard-core lemma implies regularity lemma. We have the stronger result that whatever boosting algorithm you give me, I get a hard-core lemma and regularity lemma with the same parameters and algorithm. Thus we can pick the boosting algorithm that gives the best results for our application.

Boosting with oblivious decision trees.

Let

- $f$  be function we want to learn,  $f : \chi \rightarrow \{0, 1\}$ .
- $\mathcal{T}$  class of Boolean functions we use as basic hypotheses.
- $F_k \mathcal{T}$  compositions of  $\leq k$  functions in  $\mathcal{T}$ .

Measure  $\mu : X \rightarrow [0, 1]$ ,  $|\mu| = \mathbb{E}_{x \in X} \mu(x)$ ,  $D_\mu(x) = \frac{\mu(x)}{|\mu|}$  induced distribution: pick  $x$ , keep with probability  $\mu(x)$ , else repeat.



- Weak learner: If  $\mu \in F_k \mathcal{T}$  and  $|\mu| \geq 2\delta$ , whp returns  $h \in \mathcal{J}$  so that  $\mathbb{P}_{x \sim D_\mu}[h(x) = f(x)] \geq \frac{1}{2} + \varepsilon$ .
- Strong learner: whp, finds  $H \in F_k \mathcal{T}$ ,  $\mathbb{P}_{x \in X}(H(x) = f(x)) \geq 1 - \delta$ .

Let  $\mu_0$  be constant 1.

While  $|\mu_i| \geq 2\delta$  do

- $h_{i+1} \leftarrow WL(\mu_i)$
- Partition  $X$  according to values of  $h_1, \dots, h_i$  into “blocks”  $B$ ,  $\text{Maj}_B$  most likely value in  $\delta$ .
- $H_{i+1} = \text{output } \text{Maj}_{B(x)}$ ,

$$\mu_{i+1} = \begin{cases} \frac{1-p_{\text{Maj}}}{p_{\text{Maj}}}, & \text{if } f(x) = \text{Maj} \\ 1, & \text{otherwise} \end{cases}$$

cost exponential in  $k$  in most circumstances. under what circumstances not exp in  $k$ ? If VC dimension small.

number of boolean comb. Grows at most order  $k^d$ .

Induce new distribution, make  $f$  unbiased within each block.

Total probability: on  $1 - p_{\text{Maj}}$ ,  $f = 1$ , make  $p_{\text{Maj}}$ .  $f = 0$ ,  $1 - p_{\text{Maj}}$ .

(Last time: if  $|\mu_i| \leq 2\delta$ ,  $\mathbb{P}_{x \in X}[H_i = f] \geq 1 - \delta$ .)

We need to show this method terminates.

Consider the potential function

$$\varphi = \mathbb{E}_{x \in U, B \leftarrow B(x)}[(\mathbb{P}[f = 1|B])^2]$$

Closer to 1, closer it is to boolean. This is maximized if every bucket/block is either constantly 1 or 0.

We need to show every iteration increases this potential function.

WTS: In every iteration  $\varphi$  increases by  $\Omega((\varepsilon\delta)^2)$ .

Fix  $B$ .

$$p = \mathbb{P}[f = 1|B] \quad (17)$$

$$q = \mathbb{P}[h_{i+1} = 1|B] \quad (18)$$

$$q + \alpha_+ = \mathbb{P}[h_{i+1} = 1|B, f = 1] \quad (19)$$

$$q - \alpha_- = \mathbb{P}[h_{i+1} = 1|B, f = 0] \quad (20)$$

$$\alpha_+ p = \alpha_- (1 - p) \text{ by conservation} \quad (21)$$

$$B \leftarrow \alpha_- (1 - p) \quad (22)$$

$$p_0 = \mathbb{P}[f = 1|B_0] = \frac{\mathbb{P}[f = 1 \wedge h = 0]}{\mathbb{P}[h = 0]} = \frac{p(1 - q - \alpha_+)}{1 - q} \quad (23)$$

$$p_1 = \mathbb{P}[f = 1|B_1] = \frac{\mathbb{P}[f = 1 \wedge h = 1]}{\mathbb{P}[h = 1]} = \frac{p(q + \alpha_+)}{q} \quad (24)$$

$$(\text{Value of potential}) = qp_1^2 + (1 - q)p_0^2 = p^2 \left( \frac{(q + \alpha_+)^2}{q} + \frac{(1 - q - \alpha_+)^2}{1 - q} \right) \quad (25)$$

$$= p^2 \left( \left( q + 2\alpha_+ + \frac{\alpha_+^2}{q} \right) + \left( 1 - q - 2\alpha_+ + \frac{\alpha_+^2}{1 - q} \right) \right) \quad (26)$$

$$= p^2 \left( 1 + \frac{\alpha_+^2}{q} + \frac{\alpha_+^2}{1 - q} \right) \quad (27)$$

$$\geq p^2 + 4p^2\alpha_+^2 \geq \alpha_+^2 \quad (28)$$

if  $\text{Maj} = 1$ .

Assume  $\text{Maj}_B = 1$ .

$$\mathbb{E}_{x \in B} [\mu(x)((-1)^{(h(x) \neq f(x))})] = (f = 1)p \left( \frac{1 - p}{p} \right) [(q + \alpha_+) - (1 - q - \alpha_+)] \quad (29)$$

$$+ (f = 0)(1 - p)1[1 - (1 - \alpha_-) - (q - \alpha_-)] \quad (30)$$

$$= (1 - p)[\cancel{2\eta - 1} + \cancel{1 - 2\eta} + 2\alpha_+ + 2\alpha_-] + 2\alpha_+(1 - p) + 2\alpha_+p = 2\alpha_+ \quad (31)$$

$$\mathbb{E}_{x \in B} \dots = \mathbb{E}_x \mathbb{E}_{B=B(x)} \mu(x)(-1)^{(f(x) \neq h_{i+1}(x))} \quad (32)$$

$$= \mathbb{E}_B 2\alpha_+ \geq \varepsilon |\mu| = 2\delta\varepsilon \quad (33)$$

$$\Delta\varphi = \mathbb{E}_B \alpha_t^2 \geq (\delta\varepsilon)^2. \quad (34)$$

Every step don't terminate, potential function increases by  $(\delta\varepsilon)^2$ . It's a number in  $[0, 1]$ , so the number of iterations is at most  $k \leq O\left(\frac{1}{(\delta\varepsilon)^2}\right)$ . At most polynomial in parameters.

All you get from this is a decision tree. Brute force. Complexity is exponential in  $k$ . This is a bug, not a feature.

Can always convert sign into decision tree. On hard-core, DMT give same thing. In-complexity terms, we don't get good hard-core, because circuit size  $2^k$  circuit size. Better boosting algorithm makes  $G$  have smaller complexity. If stopping point is HCL, not boosting you want. In dense model theorem, fine because all you care about is size of  $k$  not complexity of  $G$ .

This potential function matches this boosting algorithm. Other boosting algorithms you can analyze with other potential functions.

Why I want to do this way with this potential function:

- Like potential function used most in graph theory. Property: you can't make negative progress, you always go forwards.
- Translated get usual Sz, except don't get equal size partitions.

We have different batches/blocks that we're partitioning to. Suppose we get stuck: no function that correlates globally, but there are many blocks where we can find functions that correlate with the function inside that block. If  $\varepsilon$  fraction of blocks find functions that correlate, query them all. That's going to be our next round.

$2^k$  functions correlating. Now  $2^{2^k}$  buckets, but increased potential function by poly in terms of  $\varepsilon, \delta$ . This is a familiar argument, we can only go  $\frac{1}{\varepsilon}$  iterations before we terminate.

We've increased everything a tower depending on  $\varepsilon$ .  $\varepsilon$  had better be a big constant. We have a hard-core defined in terms of  $k$  things. Hard-core looks like: partition.  $1 - \delta$  fraction where have answer correct. Boolean combinations of  $k$  things you've learned. Except for  $\varepsilon$  fraction, each part: no better method than predict constant function.

Terminate: on most blocks we can't beat majority answer.

Bounds: you need a tower. No avoidance of tower in some contexts. Graphs are not the same context.

strong hardcore, or strong dense model. I don't know the formulation let alone the proof. I have the proof I need the formulation.

Take any boolean function. Any class of tests. Partition,  $\forall \varepsilon$  exists  $k$ , partition into  $k$ , such that in almost all partitions you can't do better than constant function in predicting it.

Fix set of vertices  $V$  of set  $n$ . Let  $U$  be edges in complete graph on  $V$ . (Interesting if  $U$  is not the complete graph, ex.  $U$  is the edges in  $d$ -regular expander on  $V$ .)

Let  $\tau$  be the set of cuts defined by  $A, B \subseteq V$ . WLOG  $A \cap B = \emptyset$ .

Directed edge is  $e = \{a, b\}$ .  $T(e) = 1$  iff  $a \in A$  and  $b \in B$ . Apply this chain to this class of tests.

What does  $F_K T$  look like? You have a bunch of sets  $A, B$  and you're defined by which sets they go between. Or think of it as: for every test you have 3 regions: those in  $A$ ,  $B$ , neither. When you look at all of  $k$  things, split vertices into  $3^k$  regions. A little overkill, but all the tests, can split into  $3^k$ . Knowing region of endpoints tells you which subset of tests are true.

This looks not quite equivalent. Here test sets are combinatorial rectangles. Sz: partition of left and right such that every rectangle works. That's a more structured version that what you want? One gives you the other?

If you use this kind of boosting, you don't get much better than this. If you use other types, you get more structured.

Plugging into the generic regularity lemma,  $|E| \geq \delta \binom{n}{2}$ . This says there exists  $\mu = G(T_1, \dots, T_k)$ , where  $k = O\left(\frac{1}{\varepsilon^2 \delta^2}\right)$  good predictor for any cut in the graph. This is the weak regularity of Kannan-Frieze. For Szemerédi we need the stronger boosting lemma.

Use the  $T$ 's to divide vertices into  $3^k$  subsets such that  $\mu$  is a constant on every pair of subsets. Which sets you're in dictate which is true.

$3^k$ . If you want to know how many edges between  $A$  and  $B$ , this is

$$\frac{E_G(A, B)}{|E_G|} \approx_\varepsilon \sum_{i,j} \mu_{ij} \frac{|A \cap A_i| |B \cap B_i|}{|V|^2}.$$

Constantly much information. You can't recover original set of edges, but can approximate any cut by knowing how big the intersections are relatively.

If we do this with  $G$  a subset of an expander, we get the same thing except

$$\frac{E_G(A, B)}{|E_{G'}|} =$$

Use expander mixing, with error term of expander mixing, get same thing.

### 3 Improved boosting algorithms and implications

Based on Holenstein, simplified.

We had a generic reduction from boosting algorithms to dense model theorems to regularity lemmas. Today I'll discuss how to get different intermediate results from boosting algorithms. I'll sample the variety of boosting results and see what kind of dense model theorems/regularity lemmas they give. I'll talk about recursive uses of dense model theorems to give decompositions.

I'll start with the standard algorithms and then use a modification from Holenstein necessary for the dense model theorem. The modification is fairly generic; it's about one parameter that is not important for boosting but important for dense model theorems.

A weak learner is given  $\mu_t$  and produces  $h_t$  that has some advantage on the subdistribution (normalize by expectation of the weight  $\mu_i$  to get the probability) <sup>1</sup>

$$\mathbb{E}_{x \sim_D U} (-1)^{h_t(x)=f(x)} \mu_t(x) \geq \varepsilon |\mu_i| \quad (35)$$

$$\iff \mathbb{P}_{x \sim D_{\mu_i}} [h_t(x) = f(x)] \geq \frac{1}{2} + \varepsilon \quad (36)$$

Here

$$|\mu_t| = \mathbb{E}_x |\mu_t(x)| \quad (37)$$

$$N_t(x) = \sum_{t'=1}^t (-1)^{h_{t'}(x)=f(x)} \quad (38)$$

$$I = N_t(x). \quad (39)$$

If  $I$  is negative (we're not doing so well), we weight  $x$  more; the better we're doing, the less important it is to continue doing well on the data point  $x$ .

<sup>1</sup>I think of  $U$  as uniform but it doesn't matter.

$M_I$  is a function that starts with 1, and goes to 0 gradually, for example, linear.

$$M_I = \begin{cases} 1, & \text{if } I \leq 0 \\ 1 - \gamma I, & 0 \leq I \leq \frac{1}{\gamma} \\ 0, & \text{if } I \geq \frac{1}{\gamma} \end{cases} \quad (40)$$

$$\gamma = \frac{1}{4}\varepsilon\delta. \quad (41)$$

AdaBoost works better:

$$M_I = \begin{cases} 1, & \text{if } I \leq 0 \\ (1 + \gamma)^{-I}, & I > 0 \end{cases} \quad (42)$$

$$\gamma = \frac{1}{4t}. \quad (43)$$

We need  $\gamma$  small because we need the slope to be small relative to the value.

Let  $A_I = \sum_{J \geq I} M_J$ .

$$A_t = \mathbb{E}_{x \sim U} A_{N_t(x)}$$

is our generic potential function.

$$\Delta_I = M_{I-1} - M_I \quad (44)$$

$$\Delta_t = \mathbb{E}_{x \sim U} \Delta_{N_t(x)} \quad (45)$$

How we picked  $1 + \gamma$ : to make slope proportional to value with small constant of proportionality.

$$\Delta_I = \begin{cases} 0 & \text{if } I \leq 0 \\ \gamma & \text{if } 0 \leq I \leq \frac{1}{\gamma} \\ 0 & \text{if } I > \frac{1}{\gamma} \end{cases} \quad (46)$$

$$\Delta_I = \begin{cases} 0 & \text{if } I < 0 \\ \gamma M_I & \text{if } I > 0. \end{cases} \quad (47)$$

“Move things from the negative zone to the safe region.” Potential function is how many credit we could get ourselves in the future. if nothing to give credit left to, done.

**Algorithm 3.1** (Basic boosting):     •  $\mu_0 = 1$

- Repeat until  $|\mu_t| \leq 2\delta$ .
  - $h_{t+1} \leftarrow WL(\mu_t)$
  - Increase when right, decrease when wrong:  $w_{t+1}(x) = M_{N_{t+1}(x)}$ . ( $N_t(x)$  increases or decreases by 1 depending on whether it is right at  $t$ .)
- Return  $H \in \text{Maj}(h_1, \dots, h_t)$ ,  $1 - 2\delta$ .

We'll modify this later to get  $1 - \delta$ . ( $1 - 2\delta$  is not good enough because of the 2. For learning purposes, you just replace  $\delta$  by  $\frac{\delta}{2}$ , but with the 2, it doesn't imply the hard-core set lemma or dense model theorem.)

$$N_t(x) = \sum_{t'=1}^t (-1)^{h_{t'}(x)=f(x)} \quad (48)$$

The average value is at most  $2\delta$ . Everything where majority incorrect has  $\mu$ -value 1.

We had

$$\mathbb{E}_{x \sim U} (-1)^{h_{t+1}(x)=f(x)} \mu_t(x) \geq \varepsilon |\mu_t|. \quad (49)$$

Think of this as saying: We expect most of the time to be taking a step forward rather than a step back.

If  $f(x) = h_{t+1}(x)$ , then

$$A_{N_{t+1}(x)} = A_{N_t(x)} - \mu_t(x) \quad (50)$$

If  $f(x) \neq h_{t+1}(x)$ , then

$$A_{N_{t+1}(x)} = A_{N_t(x)} + M_{N_t(x)-1} \quad (51)$$

$$= A_{N_t(x)} + M_{N_t(x)} + M_{N_t(x)} + \Delta_{N_t(x)} \quad (52)$$

For the two choices of  $M_I$ ,  $\Delta_{N_t(x)}$  is bounded by  $\gamma$ ,  $\gamma M_{N_t(x)}$ , respectively.

We get

$$\Delta A \leq -\varepsilon |\mu_t| + \mathbb{E}_x \Delta_{N_t(x)} \quad (53)$$

$$\Delta A \leq \begin{cases} -\varepsilon |\mu_t| + \gamma \\ -\varepsilon |\mu_t| + \gamma |\mu_t| \end{cases} \quad (54)$$

Using either example, we know a bound on  $\Delta$ . In example, 1, decrease of  $2\delta\varepsilon$ , so setting for  $\geq$ , for example 2,  $\delta \ll \varepsilon$  to make progress.

In the two cases

$$A_0 = \frac{1}{2\gamma} \quad t \leq O\left(\frac{1}{\varepsilon^2 \delta^2}\right) \quad (55)$$

$$A_0 = \frac{1+\gamma}{\gamma} \quad t \leq O\left(\frac{1}{\varepsilon^2 \delta}\right). \quad (56)$$

How to get rid of the factor of 2? Have a moving target for where we draw the line. As time goes on, it's harder to be in the safe zone; move things to the danger zone.

**Algorithm 3.2:**    •  $\mu_0 = 1$ ,  $s = 0$ .

- Repeat until  $|A_t| < 2\delta s$ .

$$- h_{t+1} \leftarrow WL(\mu_t)$$

- Increase when right, decrease when wrong:  $\mu_{t+1}(x) = M_{N_{t+1}(x)-s}$ .
- If  $|\mu_{t+1}| < 2\delta$ ,  $s \leftarrow s + 1$ . (Move everything 1 step to the left.)
- $\mu_{t+1}(x)M_{N_{t+1}} - s$ .
- Note that the invariant  $\mu_t > 2\delta$  is maintained in this loop.
- Return  $H_t(x)$ : if  $|\sum_{i=1}^t (-1)^{h_i}| \geq s$  use  $\text{Maj}(h_1(x), \dots, h_t(x))$  (if there are  $s$  more 1's than 0, say 1; if there are  $s$  more 0's than 1, say 0); otherwise, flip a coin with bias  $\frac{1}{2} + \frac{\sum_{i=1}^t (-1)^{h_i(s)}}{2s}$  and output 1 if heads. (Use preponderance of 0's and 1's as a coin flip.)

We saw that the change in  $A$  due to updating  $N_{t+1}$  is

$$\Delta A = \varepsilon|\mu| - \mathbb{E}\Delta_t.$$

When  $s$  is set to  $s + 1$ , the change in  $A$  is

$$\Delta A = |\mu| + \mathbb{E}\Delta_t \tag{57}$$

$$\Delta A = -\varepsilon|\mu| + |\mu| + 2\mathbb{E}\Delta_t \tag{58}$$

$$\leq 2\delta - \Omega(\varepsilon\delta) \tag{59}$$

$$\Delta(A - 2\delta s) < -\Omega(\varepsilon\delta). \tag{60}$$

This is the same decrease as before, so the bound on the step is the same as before.

Why is this good? Holenstein uses a rounding trick, also in TTB.

$$A_t(x) \geq s_t - N_t(x) + A_0, \quad \text{for } N_t(x) < s \tag{61}$$

Suppose that  $N_t(x) \geq s$ . Then we're always correct. Otherwise, we're incorrect with probability at most  $\frac{s - N_t(x)}{2s}$ . Thus

$$\mathbb{E}(\text{fraction incorrect}) = \frac{1}{2s} \mathbb{E}[s - N_t(x)] \tag{62}$$

$$\leq \mathbb{E} \frac{A_t(x)}{2s} = \frac{2\delta s}{2s} = \delta. \tag{63}$$

Is this the same as: Pick a threshold at random and use that threshold?

What did we get? We gave 2 algorithms with  $k = O(\frac{1}{\varepsilon^2\delta^2})$  and  $k = O(\frac{1}{\varepsilon^2\delta})$ . Before, we had to allow any boolean function on  $k$  variables, a large class; here, we take a majority function. This is in any complexity class that can compute majority.

$$H = g(h_1, \dots, h_k) = \text{sign}(\sum_i h_i - s) \tag{64}$$

$$\mu = G(h_1, \dots, h_k, f) = G'(\sum (-1)^{h_i(x)=f(x)}). \tag{65}$$

Both are much simpler than what we had before.

What kind of results do we get from this? Having  $g$  be in a small complexity class is important for using the hard-core set lemma for moderate complexity classes like  $TC^0$ . I'll go all the way to regularity and see what we get.

Before, we got a model of a dense graph  $|E_G| \geq \delta|V|^2$ . We have a model that looks like a partition into  $3^k$  sets. Think of this as approximating the adjacency matrix of a graph by a matrix with big blocks of constants,  $A_G = M$ ,  $\text{rank}(M) = 3^k = 3^{O(\frac{1}{\varepsilon^2\delta^2})}$ . We can replace  $\delta^2$  by  $\delta$ . We get a bunch of cuts in a graph, not necessarily disjoint, just  $\text{poly}(\frac{1}{\varepsilon^2\delta^2})$ . Weigh edges by how many of the cuts they're in. Each cut is a rank 1 matrix. If you want to know the number of cuts you're in, write as a matrix of rank  $k$  rather than  $3^k$ . The one thing we have to do: that's the crude score, we're actually interested in a piecewise linear function or exponentially decreasing function of that score. Componentwise we have to apply this function to  $M$ .

Frieze and Kannan look at a spectral approximation. Instead of boolean tests, they use tests defined by prob distributions.  $b^T A_G a$ . They get a similar result without a truncation operation.

Different boosting have different combining functions which give different structures of model which give qualitatively and quantitatively different regularity lemmas.

## 4 2017-3-2

We look at applications of dense model theorem, recursively to get a decomposition. What happens when the hypothesis to the dense model theorem fails.

Example 1.

- Let  $\mathcal{T}$  be the set of half-planes.
- Let  $S$  be a distribution on points in the unit square.
- $\delta$  dense means: For any polygon  $P$ ,  $\mathbb{P}[x \in P] \leq \frac{1}{\delta} \text{Area}(p)$ .  
(If a region is puny, then the probability of landing inside is puny.)
- $(\varepsilon, \delta)$ -pseudodense for  $\mathcal{F}_K \mathcal{T}$  means: Any region of polygons with  $\leq k$  lines and  $\mathbb{P}_{x \sim S}[x \in P] \geq \varepsilon/\delta$ .

Example 2.

- $S$  is a distribution on edges in the complete graph, say uniform on  $E(G)$ .
- $\mathcal{T}$ :  $T_{A,B}(e) = 1$  if  $e$  goes between disjoint sets  $A, B$ .
- $\delta$ -dense:  $|E(G)| \geq \delta|V|^2$ .
- $(\delta, \varepsilon)$ -pseudo-dense:

$$\frac{E_G(A, B)}{|E(G)|} \leq \frac{1}{\delta} \frac{|A||B|}{\binom{n}{2}}.$$

Example 3.

- $S$  is a distribution on  $\{0, 1\}^n$ .
- $\mathcal{T}$  are all tests that depend on  $\leq l$  variables.



- $\mathcal{F}_k\mathcal{T}$  is a subset of tests that depend on  $\leq kl$  variables.
- $\delta$ -dense means  $|S| \geq \delta 2^n$ .
- $P$  is  $(\delta, \varepsilon)$ -pseudodense if any property  $P$  that depends on  $\leq kl$  variables and has probability  $> \frac{\varepsilon}{\delta}$  in  $S$ ,  $\mathbb{P}_S(P) \leq \frac{1}{\delta} \mathbb{P}_U(P)$ .

We can learn to model to  $S$ . Is the model for  $S$  close to  $S$  or is it just a model?

The algorithmic version of the dense model theorem that we got from boosting works as follows. Let  $S$  be any distribution on  $U$ . Then for some  $k = \text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta})$ , either

1. there exists  $T \in \mathcal{F}_k\mathcal{T}$  so that

$$\mathbb{P}_{x \sim U}[T[x]] \leq \delta \mathbb{P}_{x \sim S}[T[x]] - \varepsilon$$

or

2. there exists a model  $\mu$ ,  $|\mu| \geq \delta$  (a subdistribution of the uniform distribution) and  $\mu \in \mathcal{F}_k\mathcal{T}$ ,

$$|\mathbb{P}_{x \sim S}[T(x)] - \mathbb{P}_{x \sim D_\mu}[T(x)]| \leq O\left(\frac{\varepsilon}{\delta}\right).$$

(Recall that  $D_\mu$  means: pick an element  $x$  for  $U$  and keep with probability  $\mu(x)$ .  $|\mu| = \mathbb{E}_{x \sim U} \mu(x)$ .)

To make an algorithmic version, the subroutine we need that plays the role of a weak learner is: given two distributions we can find an approximately optimal distinguisher in  $\mathcal{T}$ .

In the unit square example, the algorithmic version is pretty much guaranteed. The line for the optimal distinguisher passes through 2 of the points.

A description is a probability distribution on regions, uniform on each region. If the algorithm returns a hot spot, then we don't have a model for what's out there and what's inside. Maybe there's still a model for the outside. Use this to get a decomposition into any distribution. We want a more agnostic version, get some combination of these two things.

First, we get a decomposition lemma.

**Lemma 4.1.** *Let  $S$  be any distribution on  $U$ . Then for  $k = \text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta})$ , there exists  $T \in \mathcal{F}_K\mathcal{T}$  and  $\mu \in \mathcal{F}_k\mathcal{T}$  so that*

•

$$S \sim_{(\frac{\varepsilon}{\delta}, \mathcal{T})} D(\alpha, T, \mu)$$

, where  $\alpha = \mathbb{P}_{x \in S}[T(x) = 1]$  and  $D(\alpha, T, \mu)$  is the distribution: *ith probability  $\alpha$ , pick from  $s$  conditioned on  $T$ ; with probability  $1 - \alpha$ , pick from  $D_\mu$  conditioned on  $T$ ,*

•

$$\mathbb{P}_{x \in U} T(x) \leq \delta \mathbb{P}_{x \in S}[T(x)],$$

- $|\mu| \geq (1 - \alpha)\delta$ .

1. Start with  $T = \phi$ . Let  $\mu = 1$ . Let  $S' = S|_{\neg T}$  and  $U' = U|_{\neg T}$ .
2. Apply dense model theorem to  $S'$  within  $U'$ ,  $\alpha = \mathbb{P}_{x \in S}[T(x)]$ ,  $\delta' = (1 - \alpha)\delta$ . If we set  $\mu \in \mathcal{F}_K \mathcal{T}$  apply DMT to  $\delta', U', \varepsilon$  to get either  $\mu'$  or  $T$ .
3. If  $\mu$ , we output  $D(al, T, \mu)$ . Otherwise  $T \leftarrow T \cup T'$ .

Otherwise we get  $T'$ ,

$$\mathbb{P}_{x \in h'} T'(x) \leq \delta' \mathbb{P}_{x' \in S'} T'(x) - \varepsilon.$$

Set  $T \leftarrow T \cup T'$ . Note

$$\mathbb{P}_{x \in U}(T' - T) \leq \mathbb{P}_{x \in U}(T' - T) \tag{66}$$

$$\leq \delta' \mathbb{P}_{x \in S'}(T' - T) \tag{67}$$

$$= \delta'(1 - \alpha) \mathbb{P}_{x \in S}(T' - T) \tag{68}$$

$$= \delta \mathbb{P}_{x \in S}(T' - T). \tag{69}$$

Probability has grown: Increases by  $\alpha$  and is bounded by 1.

This is the proof of the decomposition lemma. We don't have the model yet.

Maybe we just get hotspots with a small amount of error, or just a distribution, or something in between. We're interested in when we get something in between.

Recursive decomposition lemma:

Let  $S$  be any distribution,  $j$  be an integer,  $k = \text{poly}(\frac{1}{\delta}, \frac{1}{\varepsilon}, j)$ ,  $T_1, \dots, T_j \in \mathcal{F}_k \mathcal{T}$ ,  $\mu_0, \dots, \mu_{j-1} \in \mathcal{F}_k \mathcal{T}$ ,  $\alpha_1, \dots, \alpha_j$ ,

$$T_1 \supseteq T_2 \supseteq \dots \supseteq T_j,$$

$$D(\vec{\alpha}, \vec{T}, \vec{\mu}) \sim_{O(\frac{j\varepsilon}{\delta}), \mathcal{T}} S,$$

with probability  $\alpha_0$  pick from  $\mu_0$  on  $\neg T_1$ , with probability  $\alpha_1$  pick from  $\mu_1$  on  $\neg T_2$ ,

$$\mathbb{P}_{x \in h}[T_j(x)] \leq \delta \mathbb{P}_{x \in S} T_j(S)$$

and  $\mathbb{P}_{x \in S} T_j(S) \geq \varepsilon$ . Say we know that  $S$  is  $\Delta$ -dense, choose  $\delta = \frac{1}{2}$ ,  $j = \log \frac{1}{\Delta}$ , and get

Get regularity lemma: if  $S$  is  $\Delta$ -dense, there exists a distribution  $\mathcal{F}_k \mathcal{T}$  that is  $\frac{\log(\frac{1}{\Delta})}{\varepsilon}$  indistinguishable vs.  $\mathcal{T}$  and  $k = O(\log(\frac{1}{\Delta}) \text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}))$ .

Regularity lemma that you get for graphs: Before we had a regularity lemma that said any  $\Delta$ -dense graph can be described by roughly equal partitioning using  $\exp(O(\frac{1}{\varepsilon^2 \Delta}))$  sets to estimate cut sizes.  $\exp(\text{poly}(\frac{1}{\varepsilon}, \frac{1}{\log(\frac{1}{\Delta})}))$ .

Downside: Partition of graph into this many subsets, a prob dist on pairs of subsets. Can have tiny cuts with high probability. If  $\Delta$  is small that's unavoidable. Dense graph, tiny region of original graph.