

Contents

1	Big picture	3
2	Setup	7
3	Boosting and the Hardcore lemma	9
4	Dense model theorem	10
5	Proof for boosting	11
6	Comments, Regularity lemmas	12

Wiki page: <https://simons-institute.github.io/pseudorandomness/groups/modelstructure.html>

A theme that cuts across many domains of computer science and mathematics is to find simple representations of complex mathematical objects such as graphs, functions, or distributions on data. These representations need to capture how the object interacts with a class of tests, and to approximately determine the outcome of these tests.

For example, in machine learning, the object might be a distribution on data points, high dimensional real vectors, and the tests might be half-spaces. The goal would be to learn a simple representation of the data that determines the probability of any half-space or possibly intersections of half spaces. In computational complexity, the object might be a Boolean function or distribution on strings, and the tests are functions of low circuit complexity. In graph theory, the object is a large graph, and the tests are the cuts in the graph; the representation should determine approximately the size of any cut. In additive combinatorics, the object might be a function or distribution over an Abelian group, and the tests might be correlations with linear functions or polynomials.

The focus of the working group is to understand the common elements that underlie results in all of these areas, to use the connections between them to make existential results algorithmic, and to then use algorithmic versions of these results for new purposes. For example, can we use boosting, a technique from supervised learning, in an unsupervised context? Can we characterize the pseudo-entropy of distributions, a concept arising in cryptography? Do the properties of dense graphs “relativize” to sub-graphs of expanders?

In particular, we’ll start from boosting, a technique in machine learning to go from weak learning to strong learning, i.e., taking an algorithm that learns a function only with a small correlation and making one that learns the function on almost all inputs. We’ll show how boosting implies a general Hardcore Distribution Lemma, showing that any function that cannot be $1 - \delta$ approximated by simple functions has a sub-distribution of size δ where it has almost no correlation with simple functions. By starting from boosting, we will be able to show a constructive version of this lemma. From the Hardcore Distribution lemma, we’ll derive the Dense Model Theorem used by Green and Tao to show arbitrarily long arithmetic progressions in the primes. Again, by starting with boosting, we get a general algorithmic version of DMT. This algorithmic version can then be used to derive a general

Weak Regularity Theorem, with that of Frieze and Kannan and analogs for sparse graphs as a special case.

Hopefully, at this point, the working group will segue from known connections to new connections, e.g., is there a strong boosting that implies strong regularity? Can algorithmic regularity lemmas be used in ML?

We won't assume any background and will develop everything from first principles using only simple calculations. Here's an optional reading list, and some papers we might refer to.

Papers with results we'll cover:

- Klivans and Servedio, Boosting and Hard-core Sets, FOCS 99.
- Omer Reingold, Luca Trevisan, Madhur Tulsiani, Salil P. Vadhan: Dense Subsets of Pseudorandom Sets. FOCS 2008: 76-85
- Luca Trevisan, Madhur Tulsiani, Salil P. Vadhan: Regularity, Boosting, and Efficiently Simulating Every High-Entropy Distribution. IEEE Conference on Computational Complexity 2009: 126-136
- Russell Impagliazzo, Algorithmic Dense Model Theorems and Weak Regularity
- Sita Gakkhar Russell Impagliazzo Valentine Kabanets. Hardcore Measures, Dense Models and Low Complexity Approximations

Bibliography:

We won't go through these papers explicitly, but they provide the context.

- Robert E. Schapire: The Strength of Weak Learnability (Extended Abstract). FOCS 1989: 28-33 : 01 June 2005 A decision-theoretic generalization of on-line learning and an application to boosting Yoav Freund, Robert E. Schapire
- Yoav Freund, Robert E. Schapire: Game Theory, On-Line Prediction and Boosting. COLT 1996: 325-332
- Russell Impagliazzo: Hard-Core Distributions for Somewhat Hard Problems. FOCS 1995: 538-545
- Thomas Holenstein: Key agreement from weak bit agreement. STOC 2005: 664-673
- Boaz Barak, Ronen Shaltiel, Avi Wigderson: Computational Analogues of Entropy. RANDOM-APPROX 2003: 200-215
- Alan M. Frieze, Ravi Kannan: The Regularity Lemma and Approximation Schemes for Dense Problems. FOCS 1996: 12-20
- Noga Alon, Amin Coja-Oghlan, Hiêp Hàn, Mihyun Kang, Vojtech R'odl, Mathias Schacht: Quasi-Randomness and Algorithmic Regularity for Graphs with General Degree Distributions. SIAM J. Comput. 39(6): 2336-2362(2010)

- Noga Alon, Assaf Naor: Approximating the Cut-Norm via Grothendieck’s Inequality. SIAM J. Comput. 35(4): 787-803 (2006)
- Green, Ben; Tao, Terence (2008). “The primes contain arbitrarily long arithmetic progressions”. Annals of Mathematics. 167 (2): 481547.
- Tao, Terence; Ziegler, Tamar (2008). “The primes contain arbitrarily long polynomial progressions”. Acta Mathematica. 201 (2): 213305

1 Big picture

We’ll talk about several results which have different names in different fields. You probably know them, but don’t know the same or related idea comes up in the other 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 δ -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	”	

We will take these theorems that we know to be true and show implications between them. 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

What can we gain from looking at these connections?

1. Versatility: We can “retrofit” algorithms for one setting to get algorithms for the other settings.

For example, there are many boosting algorithms. When you follow this progression, you get different quantitative and qualitative versions of dense model theorem and regularity.

2. Algorithmic and constructive results:

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

Note that the algorithmic result that we care about is different in the different settings. 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, so that we get a model that succinctly describes the set.

We pay attention to do the reductions in an algorithmic, not just an existential way.

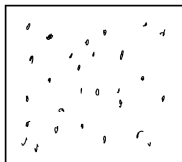
3. Using the dense model theorem for learning. Can we take a boosting technique and use it in an unsupervised way?
4. Generality: some things seem to be specific to a setting (density of graphs).

But 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. Likewise if you want to look at spectral norms rather than cuts.

Here is a cartoon:

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

Let S be some distribution on points in X .

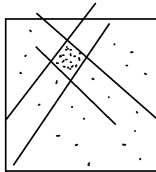


Let \mathcal{T} be a set of classifiers, ex. a set of half-planes.

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 $\mathcal{F}_K \mathcal{T}$.

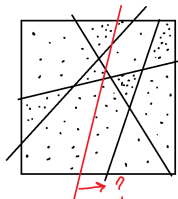
2. A violation of pseudo-density in this setting means there is a polygonal region with many more points from S than its volume, a “hot spot”.



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

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	
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	

Some comments:

1. Boosting: Note it's important that the δ 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 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.

2 Setup

First we discuss the PAC learning model.

Let U be a set, and by abuse of notation, also a distribution on that set. (Think of U as the universe, the set of possible inputs.) For simplicity, take the distribution to be uniform. Let $f : U \rightarrow \{0, 1\}$ be a boolean function. A learning algorithm can request any number of points $(x, f(x))$ where $x \sim U$. The goal is to find a hypothesis h such that

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

Theorem 2.1. *A **strong learner** for (U, f) with hypothesis class \mathcal{H} is an algorithm such that given samples $(x, f(x)), x \sim U$, outputs $h \in \mathcal{H}$ (with high probability) such that*

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

(Typically, we say that the probability of success is $1 - \varepsilon$, ask for a strong learner for all $f \in \mathcal{F}$, and require it to run in time $\text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta})$.)

In boosting, we assume that we have weak learners.

Theorem 2.2. A ε -**weak learner** for (μ, f) with hypothesis class \mathcal{H} is an algorithm such that given $(x, f(x)), x \sim \mu$, outputs h (with high probability) such that

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

It only has to output a function that is somewhat correlated with the right answer. Typically, we ask the weak learner to work on any distribution μ satisfying some assumptions.

In order to use a weak learner, we construct a routine that subsamples the distribution U to pass to the weak learner.

Definition 2.3: Let $\mu : U \rightarrow [0, 1]$. Define the probability distribution

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

1

Think of this as rejection sampling: pick $x \sim U$, keep it with probability in $[0, 1]$, or else throw it back and repeat.

In order for this sampling to be efficient, we need μ to not be too small.

Definition 2.4: Define the **density** of μ in U to be

$$|\mu| = \mathbb{E}_{x \in U} \mu(x).$$

We will use weak learners in the following context.

1. We will only run weak learners on distributions whose density is not too small (the dependence on δ is $|\mu| = \Omega(\delta)$). We don't want to run a weak learner on a distribution of very low density, because the time to simulate the distribution is inversely proportional to the density.
2. We ask the weak learners to output a function in a given class $h \in \mathcal{T}$.

Then it will turn out that both the measures that we run the weak learners on, and the final hypothesis, will be describable using $\mathcal{F}_l \mathcal{T}$ (see below), for some class \mathcal{F} .

Definition 2.5: Say that a set \mathcal{T} of functions $U \rightarrow \{0, 1\}$ form a class if $f \in \mathcal{T}$ implies $1 - f \in \mathcal{T}$.

Let \mathcal{F} be a class of boolean functions. Define the class of functions

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

¹When U is not uniform and has distribution $u(x)$, this is $\frac{\mu(x)u(x)}{\sum_{x' \in U} \mu(x')u(x')}$.

3 Boosting and the Hard-core lemma

The first boosting algorithm we 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.

We will take \mathcal{F} to be the set of all boolean functions, so given hypotheses h_1, \dots, h_k , we can choose the best predictor using $h_1(x), \dots, h_k(x)$.

Theorem 3.1 (Boosting with decision trees). *Let U be a distribution, \mathcal{T} a class of boolean functions $U \rightarrow \{0, 1\}$, \mathcal{F} the class of all boolean functions. Let $f : U \rightarrow \{0, 1\}$ be a given function (which we are trying to learn).*

1. *Suppose that there is a δ -weak learner such that given any distribution μ on U with $|\mu| \geq 2\delta$, it produces $h \in \mathcal{T}$ such that*

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

2. *Then there is a strong learner that produces $h \in \mathcal{F}_k \mathcal{T}$ with $k \leq \lceil \frac{1}{\varepsilon^2 \delta^2} \rceil$ such that*

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

2

Theorem 3.2 (Hard-core lemma). *Let U be a distribution, \mathcal{T} a class of boolean functions $U \rightarrow \{0, 1\}$, \mathcal{F} the class of all boolean functions.*

Then either

1. *There exists $h \in \mathcal{F}_k \mathcal{T}$ such that*

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

where $k \leq \frac{1}{\varepsilon^2 \delta^2}$, or

2. *(There exists a hard-core distribution.) There exists $|\mu| \geq 2\delta$ on U , such that for all $h \in \mathcal{T}$,*

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

Note it is important for us to keep track of the size of the hardcore distribution, which is $\geq 2\delta$ here. Different boosting algorithms will give the result for different classes of functions \mathcal{F} .

Proof of hard-core lemma 3.2 from boosting 3.1. Let weak learner be exhaustive search over \mathcal{T} . The weak learner operates on distributions $|\mu_i| \geq 2\delta$. If it always produces h_i with bias $\geq \delta$, then continue and obtain the strong learner: we get some $H \in \mathcal{F}_k \mathcal{T}$ such that $H(x) = f(x)$ with probability $1 - \delta$.

If at some step i our exhaustive search algorithm gets stuck, we get a distribution μ_i that's hard-core. \square

² We ignore sample complexity here. In reality, because we only see U from samples, we need to think about generalization. If the VC-dimension of \mathcal{T} is d , then the VC-dimension of $\mathcal{F}_k \mathcal{T}$ is at most k^d . In ML we don't want to take \mathcal{F} to be the class of all boolean functions. For this theorem, let's just assume we are actually given all pairs $(x, f(x))$.

4 Dense model theorem

Definition 4.1: For a set $S \subseteq U$ and a function $T : U \rightarrow \{0, 1\}$, let $T(S) := \mathbb{E}_{x \in S} T(x)$. (For a measure $\mu : U \rightarrow [0, 1]$, also write $T(\mu) = \mathbb{E}_{x \sim \mu} T(x)$.)

Let $S \subseteq U$ be a subset, and let \mathcal{T} be a set of tests. S is (ε, δ) -**pseudo-dense against** \mathcal{T} if for all $T \in \mathcal{T}$,

$$T(U) \geq \delta T(S) - \varepsilon.$$

Think of saying that the tests \mathcal{T} don't reveal that the set S is small.

1. One way of being pseudo-dense is to actually be dense.
2. Another, one step removed, is that there's a set R (or more generally, a measure μ) that's indistinguishable from S by \mathcal{T} , and such that R occupies at least a δ fraction of U .

Definition 4.2: For two distributions μ_1, μ_2 on U , we say that μ_1, μ_2 are indistinguishable by tests in \mathcal{T} up to ε , written $\mu_1 \sim_{\mathcal{T}} \mu_2$ within ε , if for every $T \in \mathcal{T}$,

$$|\mathbb{E}_{\mu_1} T - \mathbb{E}_{\mu_2} T| \leq \varepsilon.$$

Theorem 4.3 (Dense model theorem). *Let \mathcal{T} be a class of tests $U \rightarrow \{0, 1\}$.*

If S is (ε, δ) -pseudodense against $\mathcal{F}_k \mathcal{T}$, $k = O(\frac{1}{\varepsilon^2 \delta^2})$ then there exists $\mu, \mu \in \mathcal{F}_k \mathcal{T}$ such that $|\mu| \geq \frac{\delta}{1+\delta} - O(\varepsilon)$ and $D_{\mu} \sim_{\mathcal{T}} S$ to within $O(\varepsilon/\delta)$.

The idea in the proof is to use the Hard-core lemma, with the hard function being membership in S .

Proof. Let U' be the following distribution: let $\delta' = \frac{\delta}{1+\delta}$ and

1. with probability δ' , take $x \in S$ and output $(0, x)$
2. with probability $1 - \delta'$, take $x \in U$ and output $(1, x)$.

Define a test $T \in \mathcal{T}$ to operate on an example (y, x) by $T(y, x) = T(x)$. For $T \in \mathcal{F}_k \mathcal{T}$,

$$\begin{aligned} \mathbb{P}_{U'}[T((y, x)) = y] &= \delta' T(S) + (1 - \delta')(1 - T(U)) = 1 - \delta' + \delta'(T(S)) - (1 - \delta')T(U) \quad (1) \\ &= 1 - \delta' + \frac{1}{1 + \delta}(\delta T(S) - T(U)) \leq 1 - \delta' + \varepsilon. \quad (2) \end{aligned}$$

No test in $\mathcal{F}_k \mathcal{T}$ can be correct with probability $> \delta' - \varepsilon$. By the Hard-core Lemma 3.2, there exists $|\mu'| \geq 2(\delta' - \varepsilon)$ such that for any $T \in \mathcal{T}$, $\mathbb{P}_{(x,y) \sim U'}[T(x) = y] \leq \frac{1}{2} + \varepsilon$.

In order for μ' to be hardcore, it must be split approximately evenly between U and S (up to ε); otherwise; we could have an advantage by predicting constant 0 or 1. Thus each part has at least $2(\delta' - \varepsilon)(\frac{1}{2} - \varepsilon) = \delta'(1 - O(\frac{\varepsilon}{\delta}))$ of the mass. Then

$$D_{\mu'|U} \sim_{O(\varepsilon)} D_{\mu'|S} \sim_{O(\frac{\varepsilon}{\delta})} S.$$

□

5 Proof for boosting

Proof of Theorem 3.1. The algorithm is as follows. Let $WL(\mu)$ denote the weak learner operating on (μ, f) .

Let μ_0 be constant 1, $i = 0$.

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

- $h_{i+1} \leftarrow WL(\mu_i)$.
- Partition U according to values of h_1, \dots, h_i .

Let $h_{1:i}(x) := (h_1(x), \dots, h_i(x)) \in \{0, 1\}^i$, and let $B_i(x)$ be the “block” that x is in,

$$B_i(x) = h_{1:i}^{-1}(h_{1:i}(x)) = \{y \in U : h_{1:i}(x) = h_{1:i}(y)\}.$$

For a set B , let $\text{Maj}(B)$ denote the majority value of f on B .

- Define μ_{i+1} by

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

where $p_{\text{Maj}, B} = \mathbb{P}(f(y) = \text{Maj}(B) | y \in B)$, the proportion of the majority in B .

- $i \leftarrow i + 1$.

Finally, return $H_i(x) = \text{Maj}(B_i(x))$, i.e., look at the block that x is in, and choose the majority value.

Note that the measure μ_{i+1} rebalances each block B_i such that conditioned on y being in a block $B_i(x)$,

$$\mathbb{P}_{y \sim \mu_{i+1}}(f(y) = 1 | y \in B_i(x)) = \mathbb{P}_{y \sim \mu_{i+1}}(f(y) = 0 | y \in B_i(x)) = \frac{1}{2}.$$

Indeed, we have

$$\mathbb{E}_{y \sim U} [\mathbb{1}_{f(y)=1} \mu_{i+1}(y) | y \in B_i(x)] = p_{\text{Maj}, B_i(x)} \frac{1 - p_{\text{Maj}, B_i(x)}}{p_{\text{Maj}, B_i(x)}} = 1 - p_{\text{Maj}, B_i(x)} \quad (3)$$

$$\mathbb{E}_{y \sim U} [\mathbb{1}_{f(y)=0} \mu_{i+1}(y) | y \in B_i(x)] = (1 - p_{\text{Maj}, B_i(x)}) \cdot 1 = 1 - p_{\text{Maj}, B_i(x)} \quad (4)$$

$$|\mu_{i+1}| = \mathbb{E}_{y \sim U} [\mu_{i+1}(y)] = \sum_{\text{block } B_i} [2(1 - p_{\text{Maj}, B_i}) \mathbb{P}(B_i)] \quad (5)$$

$$\geq 2(1 - p_{\text{Maj}, U}). \quad (6)$$

Note that if $|\mu_{i+1}| \leq 2\delta$, then $\mathbb{P}_{x \in X}[H_i = f] \geq 1 - \delta$, and we are done. (We stop before we have to apply the weak learner to a distribution of density $< \delta$.)

We need to show this method terminates in a bounded number of steps.

Consider the potential function

$$\varphi_i = \mathbb{E}_{x \sim U} [(\mathbb{P}[f = 1 | B_i(x)])^2] = \mathbb{E}_{x \sim U} [\mathbb{E}[f | B_i]^2]$$

(Think of B_i as a partition; for a partition, $\mathbb{E}[f|P]$ is a function of x that takes x to the average value in the atom of the partition that contains x .) Note this have value in $[0, 1]$ and is maximized if f is constant on every block. We show every iteration increases this potential function by at least a fixed amount, $(\varepsilon\delta)^2$.

Fix a block B in the partition. Define $p, q, \alpha_+, \alpha_-, p_0, p_1$ as follows.

$$p = \mathbb{P}[f = 1|B] \tag{7}$$

$$q = \mathbb{P}[h_{i+1} = 1|B] \tag{8}$$

$$q + \alpha_+ = \mathbb{P}[h_{i+1} = 1|B, f = 1] \tag{9}$$

$$q - \alpha_- = \mathbb{P}[h_{i+1} = 1|B, f = 0] \tag{10}$$

$$\alpha_+p = \alpha_-(1 - p) \text{ by conservation} \tag{11}$$

$$p_0 = \mathbb{P}[f = 1|h = 0, B] = \frac{\mathbb{P}[f = 1 \wedge h = 0|B]}{\mathbb{P}[h = 0|B]} = \frac{p(1 - q - \alpha_+)}{1 - q} \tag{12}$$

$$p_1 = \mathbb{P}[f = 1|h = 1, B] = \frac{\mathbb{P}[f = 1 \wedge h = 1|B]}{\mathbb{P}[h = 1|B]} = \frac{p(q + \alpha_+)}{q} \tag{13}$$

$$\mathbb{E}_{x \in B}[\mathbb{E}[f|B_{i+1}]^2] = qp_1^2 + (1 - q)p_0^2 = p^2 \left(\frac{(q + \alpha_+)^2}{q} + \frac{(1 - q - \alpha_+)^2}{1 - q} \right) \tag{14}$$

$$= p^2 \left(\left(q + 2\alpha_+ + \frac{\alpha_+^2}{q} \right) + \left(1 - q - 2\alpha_+ + \frac{\alpha_+^2}{1 - q} \right) \right) \tag{15}$$

$$= p^2 \left(1 + \frac{\alpha_+^2}{q} + \frac{\alpha_+^2}{1 - q} \right) \tag{16}$$

$$\geq p^2 + 4p^2\alpha_+^2 \geq p^2 + \alpha_+^2 \tag{17}$$

$$\mathbb{E}[f|B_{i+1}]^2 - \mathbb{E}[f|B_i]^2 = \alpha_+^2(B_i(x)). \tag{18}$$

Assume WLOG that $\text{Maj}(B_i(x)) = 1$. (Otherwise the LHS is smaller.)

$$\mathbb{E}_{x \in B} [\mu(x)((-1)^{h(x) \neq f(x)})] = p \left(\frac{1 - p}{p} \right) [(q + \alpha_+) - (1 - q - \alpha_+)] \quad (f = 1) \tag{19}$$

$$+ (1 - p)1[1 - (1 - \alpha_-) - (q - \alpha_-)] \quad (f = 0) \tag{20}$$

$$= (1 - p)(2\alpha_+ + 2\alpha_-) \tag{21}$$

$$= 2\alpha_+(1 - p) + 2\alpha_+p = 2\alpha_+ \tag{22}$$

$$\mathbb{E}_{x \sim U} 2\alpha_+(B_i(x)) = \mathbb{E}_{x \sim U} [\mu(x)((-1)^{h(x) \neq f(x)})] \tag{23}$$

$$\geq \varepsilon|\mu| \geq 2\delta\varepsilon \tag{24}$$

$$\varphi_{i+1} - \varphi_i \geq \mathbb{E}_{x \sim U} [\mathbb{E}[f|B_{i+1}]^2 - \mathbb{E}[f|B_i]^2] \tag{25}$$

$$\geq \mathbb{E}_{x \sim U} \alpha_+^2(B_i(x)) \geq (\delta\varepsilon)^2. \tag{26}$$

Because φ_i is always in $[0, 1]$, the number of iterations is at most $k \leq (\delta\varepsilon)^2$. \square

6 Comments, Regularity lemmas

Some comments:

1. All you get from this proof is a decision tree; the complexity is exponential in k . This is a bug, not a feature.

In complexity terms, we don't get good hard-core measure, because the circuit size for the outer function G is 2^k . A better boosting algorithm would give G have smaller complexity. If your stopping point is the hard-core lemma, this is not the boosting algorithm you want. For the dense model theorem, this is fine because all you care about is size of k , not the complexity of G .

There is another boosting algorithm which gives a weighted majority function, which is a simpler function. A weighted majority can be converted into a decision tree, but not vice versa.

2. This potential function matches this boosting algorithm. Other boosting algorithms can be analyzed with other potential functions.

This is like the potential function used most in graph theory. Key property: you can't make negative progress; you always go forwards.

3. For Szemerédi regularity, we need a stronger boosting theorem. Suppose we get stuck at some step: no function correlates globally, but there are many blocks where we can find functions that correlate with the function inside that block. If in ε fraction of blocks we find functions that correlate, partition them based on all the values of these functions, and repeat.

In one step we've gone from order of 2^k to order of 2^{2^k} buckets, and increased the potential function by a polynomial in terms of ε, δ . This is a familiar argument; we can only go $\frac{1}{\varepsilon}$ iterations before we terminate. This time, the number of sets is a tower depending on ε .

4. Regularity lemmas:

Fix a set of vertices V of set n . Let U be edges in complete graph on V . (We can also consider the case when U is not the complete graph, ex. U is the edges in d -regular expander on V .)

The underlying set we care about is the set of cuts defined by $A, B \subseteq V$ where $A \cap B = \emptyset$; there are 3^k of them.

If $|E| \geq \delta \binom{n}{2}$, the generic regularity lemma says there exists $\mu = G(T_1, \dots, T_k)$, where $k = O\left(\frac{1}{\varepsilon^2 \delta^2}\right)$, that is a good predictor the number of edges of any cut in the graph. Use the T 's to divide the vertices into 3^k subsets such that μ is a constant on every pair of subsets.

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

This is the weak regularity of Frieze-Kannan. For Szemerédi we need the stronger boosting lemma (see previous point).

We can also do something similar with G a subset of an expander. The expander mixing lemma gives an error term.