

AN OPTIMAL “AIN’T OVER TILL IT’S OVER” THEOREM

RONEN ELDAN*, AVI WIGDERSON†, AND PEI WU‡

ABSTRACT. We study the probability of Boolean functions with small max influence to become constant under random restrictions. Let f be a Boolean function such that the variance of f is $\Omega(1)$ and all its individual influences are bounded by τ . We show that when restricting all but a $\rho = \tilde{\Omega}((\log 1/\tau)^{-1})$ fraction of the coordinates, the restricted function remains nonconstant with overwhelming probability. This bound is essentially optimal, as witnessed by the tribes function $\text{TRIBE} = \text{AND}_{n/C \log n} \circ \text{OR}_{C \log n}$.

We extend it to an anti-concentration result, showing that the restricted function has nontrivial variance with probability $1 - o(1)$. This gives a sharp version of the “it ain’t over till it’s over” theorem due to Mossel, O’Donnell, and Oleszkiewicz. Our proof is discrete, and avoids the use of the invariance principle.

We also show two consequences of our above result: (i) As a corollary, we prove that for a uniformly random input x , the block sensitivity of f at x is $\Omega(\log 1/\tau)$ with probability $1 - o(1)$. This should be compared with the implication of Kahn, Kalai and Linial’s result, which implies that the average block sensitivity of f is $\Omega(\log 1/\tau)$. (ii) Combining our proof with a well-known result due to O’Donnell, Saks, Schramm and Servedio, one can also conclude that: Let $\rho = \tilde{\Omega}(1/\sqrt{\log(1/\tau)})$. Restricting all but a ρ fraction of the coordinates of a monotone function f , then the restricted function has decision tree complexity $\Omega(\tau^{-\Theta(\rho)})$ with probability $\Omega(1)$.

*Microsoft Research. Partially supported by NSF grant CCF-1900460.
E-mail: roneneldan@microsoft.com.

†Institute for Advanced Study, Princeton, NJ 08540. Supported by NSF grant CCF-1900460.
E-mail: avi@ias.edu.

‡Institute for Advanced Study, Princeton, NJ 08540. Supported by NSF grant CCF-1900460.
E-mail: pwu@ias.edu.

1. INTRODUCTION

For any Boolean function $f : \{-1, 1\}^n \rightarrow \{0, 1\}$, the individual *influence* of the i th coordinate is the probability of flipping the value of f by flipping x_i on a random input x . Let $x \oplus (-1)^{e_i}$ denote the string obtained by flipping the i th coordinate of x , then

$$I_i(f) := \mathbb{P}_{x \in \{-1, 1\}^n} [f(x) \neq f(x \oplus (-1)^{e_i})].$$

In this paper, we study Boolean functions with small influences, hence functions satisfying

$$I_\infty(f) := \max_{i \in [n]} I_i(f) = o(1).^1$$

Let \mathcal{R}_p denote a p -random restriction, namely, a randomly-chosen subcube where for each coordinate, one flips a coin and, with probability p one fixes the value of the coordinate to -1 or 1 (with equal probabilities) and, with probability $1 - p$, the coordinate is left undetermined (alive). Then $f|_{\mathcal{R}_p}$ is a random sub-function given by restricting f to the subcube.

In this paper, we study Boolean functions with small influences under random restrictions. Our main goal is to prove a lower bound for the probability of the function to remain nonconstant under the restriction. We prove the following near-optimal result:

Theorem 1.1 (A simplified version of Theorem 4.1). *Given $f : \{-1, 1\}^n \rightarrow \{0, 1\}$ such that the variance of f is $\Omega(1)$, and $\tau := I_\infty(f) = o(1)$. Let $\mathcal{R}_{1-\rho}$ be a random restriction where*

$$\rho = \Omega\left(\frac{\log \log(1/\tau)}{\log(1/\tau)}\right).$$

Then for any $p \geq I_\infty(f)^{\Theta(\rho)}$,

$$\mathbb{P}[\text{Var}[f|_{\mathcal{R}_{1-\rho}}] \leq p^{\tilde{\Theta}(\frac{1}{\rho})}] \leq p.$$

The bound on the variance is near-optimal by the majority function. In particular, if f is the majority function, then

$$\mathbb{P}[\text{Var}[f|_{\mathcal{R}_{1-\rho}}] \leq p^{\Theta(\frac{1}{\rho})}] \leq p.$$

Furthermore, our bound on ρ is optimal up to a $\log \log$ factor. Because randomly restricting the tribes function with $\rho = O(1/\log(1/\tau))$, we get a constant function with probability $\Omega(1)$. Previously, Mossel et al. proved a similar result for $\rho = \Omega(1/\sqrt{\log(1/\tau)})$ using completely different techniques [17]. Prior to Mossel et al.’s work, the related conjecture, with a very suggestive name “it ain’t over till it’s over” conjecture, was proposed by Kalai and Friedgut in studying social indeterminacy [13, 14]. It implies a quantitative version of the Arrow’s Theorem. We refer the interested readers to [17] for more discussions.

Next, we discuss a corollary of this theorem to block sensitivity of functions with small influences. The *sensitivity* of an input x with respect to Boolean function f ,

¹For the rest of the paper, we consider the function f as a family of functions. Thus here by $o(\cdot)$, we mean “as n goes to infinity.” The bound $o(1)$ on the influences is worse than needed. For illustration, this is good enough as many examples we are interested in this paper satisfy that their influences are $o(1)$.

denoted $s_f(x) := \sum_{i \in [n]} [f(x) \neq f(x \oplus (-1)^{e_i})]$, is the number of the of Hamming neighbors of x which have a different function value. An inequality by Kahn, Kalai and Linial [12] asserts that

$$\mathbb{E}_x[s_f(x)] = \Omega\left(\log \frac{\text{Var}(f)}{I_\infty(f)}\right),$$

which naturally leads to the question of whether it is also true that $s_f(x) = \Omega\left(\log \frac{\text{Var}(f)}{I_\infty(f)}\right)$ for a *typical* point x . This is clearly not the case, as witnessed by the majority function. However, a corollary to our theorem is that such an estimate does indeed hold true for most points x , if sensitivity is replaced by the related notion of *block sensitivity*.

The block sensitivity of an input x with respect to function f , denoted $\text{bs}_f(x)$ is the maximum number of disjoint sets $S_1, S_2, \dots, S_m \subseteq [n]$, such that for $i \in [m]$, one has $f(x) \neq f(x \oplus (-1)^{1_{S_i}})$, by $x \oplus (-1)^{1_{S_i}}$ we mean flip the sign of variables in S_i . Clearly,² one has $\text{bs}_f(x) \geq s_f(x)$ for all f, x . Our second result shows that for functions with small influences, the block sensitivity is large on almost all points x :

Theorem 1.2. *For any function $f : \{-1, 1\}^n \rightarrow \{0, 1\}$ such that its variance is $\Omega(1)$, and $\tau := I_\infty(f) = o(1)$. Then*

$$\mathbb{P}_x[\text{bs}_f(x) \geq \tilde{\Omega}(\log 1/\tau)] = 1 - o(1).$$

Finally, if the function f is monotone in addition to having small influences, our analysis to Theorem 1.1 implies an upper bound on the influences of f under random restrictions. In the work due to O’Donnell et al. [18], it is proved that every shallow decision tree must have an influential variable. Combining these facts, one can also conclude that, for monotone function f , the restricted function will have large decision tree complexity. In particular, let $\text{DT}(f)$ denote the decision tree complexity of f . Then,

Theorem 1.1. *For any monotone function $f : \{-1, 1\}^n \rightarrow \{0, 1\}$ with $\Omega(1)$ variance, and $\tau = I_\infty(f) = o(1)$. Then for any $\rho = \tilde{\Omega}(\sqrt{1/\log(1/\tau)})$,*

$$\mathbb{P}[\text{DT}(f|_{\mathcal{R}_{1-\rho}}) = \tau^{-\Theta(\rho)}] \geq \frac{1}{2}.$$

The above theorem is, in a sense, a reverse statement to the Håstad switching lemma, which states that applying the $(1 - O(1/\log n))$ -random restriction to any polynomial-size DNF/CNF (or in general any AC^0 circuits), one gets a shallow decision tree with high probability. Our result, on the contrary, states that random restrictions with alive probability $\tilde{\Omega}(1/\log(I_\infty(f)))$ cannot simplify f to a too shallow decision tree for monotone functions f with low influences.

Context and related works. The notion of influences studied in this paper is first introduced by Ben-Or and Linial [4] in the context of *collective coin flipping*. It coincides with the “Banzaf index” studied in game theory. The class of Boolean functions with small influences have been widely studied. There are several motivations to study such functions. First, they arise naturally in social choice theory [13, 14]. For example, in a voting system of two candidates and n voters, each bit x_i represents the individual preference of each voter between the two candidates. When aggregating the social preference, it is natural to use a function

²By taking singleton sets above.

f where the potential of any given individual to determine the final outcome is limited. Second, from an algorithmic perspective, suppose that we have access to the input via a limited number of queries. Then, it is natural to query a variable when its individual influence is large. In many cases, such variables can be found iteratively and this process leads to a good approximation of f with a small number of queries. This observation has been applied in different settings [8, 1]. In computational complexity for example, to distinguish the dictatorship function v.s. functions with small individual influences is a key component of proving optimal NP-hardness for approximations [3, 10, 11, 15]. From an analytic perspective, it has been observed that functions with small influences exhibit improved concentration inequalities (e.g., [20]) and often tend to exhibit Gaussian-like behavior [17].

Applying random restrictions and studying the properties of the restricted functions has been widely studied and has led to breakthroughs in a variety of areas. For example, it is the key idea of the exponential lower bounds in circuit complexity [9] and the dramatic improvements of the sunflower lemma in combinatorics [2].

The problem of determining whether a function with small influences becomes constant under random restrictions has attracted some attention in the context of hardness amplification within NP for circuits [19].

A sub-optimal version of Theorem 1.1 follows from the “It ain’t over till it’s over” theorem proven by Mossel, O’Donnell, and Oleszkiewicz in [17]. Their approach uses the *invariance principle*, which at a high level asserts that when feeding a “smooth”³ function f with independent random inputs X_1, X_2, \dots, X_n from a product space such that each X_i has zero mean, unit second moment and bounded third moment, then the output distribution is “invariant” to the actual distribution of the inputs. This approach usually studies a related problem, then translates the result of the related problem to the Boolean cube. This translation suffers from two drawbacks. First, it obscures what is actually happening in Boolean cube. Second, the requirement of f being “smooth” normally requires additional technical treatment, and becomes the main obstacle for obtaining an optimal result.

Our approach. Our approach relies on a control-theory point of view to the problem combined with ideas from “pathwise-analysis,” using arguments which are somewhat inspired by [7]. We assume that the coordinates are revealed in a random order and are randomly assigned values ± 1 one by one. For each coordinate being revealed, we assume that with probability p a player gets an opportunity to “override” the value that has been assigned to that coordinate.⁴ If the player has the capability of deciding, with high probability, the value of the function, this implies that restricting all but a p -fraction of coordinates leaves the restricted function nonconstant.

To this end, assume that $X(t) \in \{-1, 0, 1\}^n$ is the process where at each step another coordinate is being revealed, where coordinates whose value is not determined are set to 0. We view our Boolean function f as a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ by considering its multilinear extension. If the player does not override any coordinate,

³By “smooth” here, we mean f has low degree. With additional work, the invariance principle applies to f that has its Fourier mass concentrated in low degrees.

⁴We note that Lichtenstein-Linial-Saks [16] also study a control theoretic problem, which in the surface may seem similar. The main difference is that in their model the player picks which coordinates to influence, whereas in our model these coordinates are picked randomly, as we will see momentarily. The two models differ drastically in their nature.

then the process $M(t) := f(X(t))$ is a martingale (where, for $x \in \{-1, 0, 1\}^n$, the expression $f(x)$ is the value of f when taking expectation over coordinates whose value is set to 0). The player’s ability to override coordinates effectively allows the player to add a drift to $M(t)$, where the player’s goal is to end up with $M(n)$ being equal to 1 (or, by the same argument, to 0 simply by replacing f with $1 - f$).

At this point, let us assume for simplicity that the increments of the martingale $M(t)$ have a fixed step size η (in other words, assume that it is actually a random walk up to the time when it hits $\{0, 1\}$). Moreover assume that $M(0)$ is bounded away from $\{0, 1\}$. Suppose that the player has probability p to override each step and is trying to force the process to end up at the value 1 by overriding the increment with the value $+\eta$. In this case over t steps the process accumulates a variance of $\eta^2 t$ and a drift of size $t p \eta$. Since the process eventually moves a distance of $\Omega(1)$, and hence accumulates a variance of constant order, we have $t \asymp \eta^{-2}$. It follows that in order for the effect of the drift to be more significant than that of the variance, one arrives at the condition $p \gg \eta$. In other words, the process can be efficiently controlled (meaning that the player gets to determine its endpoint) as long as the step size is at most p . In fact, we will see that this heuristic is only correct when $M(t)$ is not close to the edges, which will create an additional technical complication.

The step size of the process is, in turn, is controlled by the ℓ_∞ norm of the first-order Fourier coefficients of restrictions of the function f , or equivalently by the quantity $\max_i |\partial_i f(X(t))|$, where i is over the coordinates not fixed at time t . We need to show that this quantity remains small along the process, which is where the fact that the initial influences are small will be used.

The control of the first-order Fourier coefficients relies on a new *hypercontractive inequality* for random restrictions. We consider random restrictions $\mathcal{R}_p, \mathcal{R}_q$, where $0 \leq p \leq q \leq 1$ are the probabilities of a variable being fixed, and show that for any multilinear function f and any $0 \leq \epsilon \leq q - p$,

$$(\mathbb{E}[\mu(f|\mathcal{R}_p)^{2+\epsilon}])^{\frac{1}{2+\epsilon}} \leq (\mathbb{E}[\mu(f|\mathcal{R}_q)^2])^{\frac{1}{2}}, \quad (1.1)$$

where we use $\mu(f)$ to denote the expected value of f over the uniform measure on $\{-1, 1\}^n$.

The hypercontractive inequality will allow us to control the evolution of the first-order coefficients under the original (namely, the uncontrolled) martingale. However, we need to control those coefficients under the modified process (where the player gets to override some coordinates). This can be solved by assuming that the strategy taken by the player tries to mimic yet another process obtained by *conditioning* the original martingale $M(t)$ to end up at the value 1 (0, respectively). This amounts to a change of measure over the space of paths of $X(t)$ which gives tractable formulas for the corresponding change of measure of a single step. Equivalently, this is the strategy which ensures ending up at the desired value under a change of measure which has the minimal possible relative entropy to the uncontrolled process.

Finally, we explain how to strengthen the above result to give a quantitative bound on the variance of the restricted function. We analyze the Kullback-Leibler divergence between, roughly speaking, the string $Y(n)$ generated by the “controlled” process given the restrictions $\mathcal{R}_{1-\rho}$ determined by those coordinates that is not controlled by the player, and a uniformly random string $X \in \{-1, 1\}^n$. With the

Fourier-analytic tool of Level-1 inequality, one can show that the expected KL-divergence over the random restrictions is about $\tilde{O}(1/\rho)$. Somewhat surprisingly, the KL-divergence is, in addition to being small in expectation, highly concentrated. Recall that $Y(n)$ is sampled from $f^{-1}(1)$. All these imply that $\mu(f|_{\mathcal{R}_{1-\rho}}) \geq \exp(-\tilde{O}(1/\rho))$ with high probability. The variance bound then follows immediately once we put together with the other direction that $\mu(f|_{\mathcal{R}_{1-\rho}}) \leq 1 - \exp(-\tilde{O}(1/\rho))$ by replacing f with $1 - f$.

Organization. We present the necessary preliminaries in Section 2. Then in Section 3, we carefully define the uncontrolled and controlled process discussed in the introduction and we study the properties of these random processes. With this tool at our disposal, we prove our main result Theorem 1.1 in Section 4. Then we explain the applications of this result to the block sensitivity and decision tree complexity. We leave the technical analysis to the final section, Section 5, that the Fourier coefficients of the first order remain small under random restrictions.

2. PRELIMINARIES

General. We adopt the shorthand notation $[n]$ for the set $\{1, 2, \dots, n\}$. For a string $x \in \{-1, 1\}^n$ and a set $S \subseteq \{1, 2, \dots, n\}$, we let $x|_S$ denote the restriction of x to the indices in S . In other words, $x|_S = x_{i_1}x_{i_2}\dots x_{i_{|S|}}$, where $i_1 < i_2 < \dots < i_{|S|}$ are the elements of S . Analogously, for any function $f : \Omega \rightarrow \mathbb{R}$ over an arbitrary domain Ω . Let $A \subseteq \Omega$, we adopt the notation $f|_A$ for the sub-function of f over the domain A . Namely, $f|_A(x) = f(x)$ for $x \in A$. Given a set S , when the universe U is clear from the context we use $\bar{S} := U \setminus S$ to denote the complement of S . The *characteristic function* of a set S is given by

$$\mathbf{1}_S(i) = \begin{cases} 1 & \text{if } i \in S, \\ 0 & \text{otherwise.} \end{cases}$$

For a permutation $\pi : U \rightarrow U$. Let πS be the permuted set of S , i.e., $\pi S = \{\pi(i) : i \in S\}$.

The primary interest of this paper is Boolean functions $f : \{-1, 1\}^n \rightarrow \{0, 1\}$. Note that we use $-1, 1$ to denote “true” and “false” on the domain of f , respectively. For example, the logic AND function and the logic OR function are defined as below,

$$\bigwedge_{i=1}^n x_i = \begin{cases} 1 & x_i = -1 \ \forall i \in [n], \\ 0 & \text{otherwise,} \end{cases} \quad \bigvee_{i=1}^n x_i = \begin{cases} 0 & x_i = 1 \ \forall i \in [n], \\ 1 & \text{otherwise.} \end{cases}$$

We abuse the notation $x \oplus y$ to denote the entrywise XOR function for $x, y \in \{-1, 1\}^n$. Thus $(x \oplus y)_i = x_i \cdot y_i$. For any univariate function $h : \mathbb{R} \rightarrow \mathbb{R}$ and $x \in \mathbb{R}^n$, the application of h to x means entrywise application, i.e., $h(x)$ denotes the vector such that $(h(x))_i = h(x_i)$. For any $f : \{-1, 1\}^n \rightarrow \mathbb{R}$, $S \subseteq \{1, 2, \dots, n\}$ and $y \in \{-1, 1\}^n$, let $\text{cube}_{S,y} := \{x \in \{-1, 1\}^n : x|_S = y|_S\}$ be the subcube of $\{-1, 1\}^n$. We abbreviate $f|_{(S,y)} = f|_{\text{cube}_{S,y}}$. The same definition $f|_{(S,y)}$ extends to $y \in \mathbb{R}^T$ for any $T \geq S$ such that $y|_S \subseteq \{-1, 1\}^S$. A random p -restriction \mathcal{R}_p is a random tuple (S, y) such that for each $i \in [n]$, $i \in S$ with probability p , and y is a uniformly random element from $\{-1, 1\}^n$.

For a logical condition C , we use the Iverson bracket

$$\mathbb{I}\{C\} = \begin{cases} 1 & \text{if } C \text{ holds,} \\ 0 & \text{otherwise.} \end{cases}$$

Denote $|x|$ the length of x for any vector $x \in \mathbb{R}^n$, i.e.,

$$|x| = \left(\sum_{i=1}^n x_i^2 \right)^{1/2}.$$

For two vectors $x, y \in \mathbb{R}^n$, we adopt the following inner product

$$\langle x, y \rangle = \sum_{i \in [n]} x_i y_i.$$

The set $\{e_1, e_2, \dots, e_n\}$ forms a standard basis, where e_i denotes the vector whose only nonzero coordinate i is 1.

Given some space Ω and a probability measure γ over Ω . If the random variable X is drawn from γ , we denote it by $X \sim \gamma$. For any function $f : \Omega \rightarrow \mathbb{R}$, we often abbreviate the expectation of f over γ as $\gamma(f)$, namely,

$$\gamma(f) := \int_{x \in \Omega} f(x) d\gamma.$$

We let $\ln x$ and $\log x$ stand for the natural logarithm of x and the logarithm of x to base 2, respectively. For any distribution γ over some discrete space Ω , the entropy function

$$H(\gamma) = \mathbb{E}_{x \in \Omega} \gamma(x) \log \frac{1}{\gamma(x)}.$$

When Ω contains only two elements, we can think of the binary entropy function $H : [0, 1] \rightarrow [0, 1]$ is given by

$$H(x) = x \log \frac{1}{x} + (1 - x) \log \frac{1}{1 - x}.$$

Basic calculus reveals that for $x \in [-1, 1]$,

$$1 - H(x) \leq 4 \left(x - \frac{1}{2} \right)^2. \quad (2.1)$$

Recall that the Kullback-Leibler divergence (KL-divergence) between two distributions μ_0, μ_1 over Ω is defined by the following formula

$$\text{KL}(\mu_0 \parallel \mu_1) = \sum_{x \in \Omega} \mu_0(x) \log \frac{\mu_0(x)}{\mu_1(x)}.$$

The KL-divergence is convex. In particular, let $\mu_0, \mu_1, \gamma_0, \gamma_1$ be distributions over the same space. Then for any $\lambda \in [0, 1]$,

$$\text{KL}(\lambda \mu_0 + (1 - \lambda) \mu_1 \parallel \lambda \gamma_0 + (1 - \lambda) \gamma_1) \leq \lambda \text{KL}(\mu_0 \parallel \gamma_0) + (1 - \lambda) \text{KL}(\mu_1 \parallel \gamma_1).$$

If two random variables X_0, X_1 obey μ_0 and μ_1 , respectively, we also use $\text{KL}(X_0 \parallel X_1)$ to denote the KL-divergence between the two distributions. The KL-divergence satisfies the following chain rule:

$$\text{KL}(X_0 Y_0 \parallel X_1 Y_1) = \text{KL}(X_0 \parallel X_1) + \mathbb{E}_{x \sim X_0} \left[\text{KL} \left(\frac{Y_0 \mid X_0 = x}{Y_1 \mid X_1 = x} \right) \right].^{5,6}$$

The following simple analytical fact will be useful for us.

Fact 2.1. *Given $x, p \in \mathbb{R}$, then*

- (i) $(1+x)^p \geq 1+xp$, for any $x > -1$, and $p \geq 1$.
- (ii) $(1+x)^p \leq 1+xp$, for any $x > -1$, and $0 \leq p \leq 1$.

Proof. Let $g = (1+x)^p - 1 - xp$ be a function on x . Then

$$g' = p(1+x)^{p-1} - p. \quad (2.2)$$

When $p > 1$, (2.2) is negative for $x \in (-1, 0)$ and nonnegative for $x \geq 0$. Thus, g is decreasing in the interval $x \in (-1, 0)$ and increasing in $(0, \infty)$. Plug $x = 0$ into g , we get 0. Therefore, $(1+x)^p \geq 1+xp$. When $0 < p < 1$, (2.2) is positive for $x \in (-1, 0)$ and nonpositive for $x \geq 0$. Hence g obtains its maximum within $(-1, \infty)$ at point $x = 0$. \square

Discrete Fourier analysis. Let $f : \{-1, 1\}^n \rightarrow \{0, 1\}$ be any Boolean function. We would often treat f as a function $f : [-1, 1]^n \rightarrow [0, 1]$ or $f : \mathbb{R}^n \rightarrow \mathbb{R}$ by considering its multilinear extension, i.e.,

$$f(x) = \sum_{S \subseteq [n]} \hat{f}(S) \chi_S,$$

here χ_S is the abbreviation of $\prod_{j \in S} x_j$. An important observation is that under this notation,

$$f(0) = \mathbb{E}_{x \in \{-1, 1\}^n} [f(x)].$$

The set $\{\chi_S\}_{S \subseteq [n]}$ is a complete orthogonal basis of the space $\mathbb{R}^{\{-1, 1\}^n}$. Further,

$$2^{-n} \langle \chi_S, \chi_T \rangle = \begin{cases} 1 & S = T, \\ 0 & S \neq T. \end{cases}$$

Thus, for any $f, g : \{-1, 1\}^n \rightarrow \mathbb{R}$, we have the following by Parseval’s identity and Plancherel Theorem,

$$2^{-n} \langle f, g \rangle = \sum_{S \subseteq [n]} \hat{f}(S) \hat{g}(S). \quad (2.3)$$

$$\mathbb{E}_{x \in \{-1, 1\}^n} [f^2] = \sum_{S \subseteq [n]} \hat{f}(S)^2. \quad (2.4)$$

Adopt the following notations for partial derivatives and the vector differential operator,

$$\begin{aligned} \partial_i f(x) &= \sum_{S \ni i} \hat{f}(S) \chi_{S \setminus \{i\}}, \\ \nabla f &= (\partial_1 f, \partial_2 f, \dots, \partial_n f). \end{aligned}$$

⁵We refer the interested readers to [6] for a complete treatment.

⁶Here, we use the fraction-like notation to also denote the KL-divergence for aesthetics, as we are comparing two conditional distributions. The numerator in the fraction-like notation corresponds to the first argument in the standard notation.

For functions on Boolean cubes, by considering their multilinear extensions it’s easy to see that the above definitions work exactly as expected: For any $\delta \in \mathbb{R}$,

$$f(x + \delta e_i) - f(x) = \delta \cdot \partial_i f(x).$$

An important fact about the weight of the Fourier coefficients is the following inequality, often referred to as the Level-1 inequality.

Theorem 2.2 (Level-1 inequality). *Let $f : [-1, 1]^n \rightarrow \{0, 1\}$ be the multilinear extension of a Boolean function. Then for some absolute constant C , we have*

$$|\nabla f(0)|^2 \leq C f(0)^2 \ln \frac{e}{f(0)}.$$

We adopt the following standard definitions of the individual influence and the max influence of function f :

$$\begin{aligned} I_i(f) &= \mathbb{E}_{x \in \{-1, 1\}^n} [\partial_i f(x)^2]. \\ I_\infty(f) &= \max_{i \in [n]} I_i(f). \end{aligned}$$

By Plancherel Theorem,

$$I_i(f) = \sum_{S \subseteq [n]: i \in S} f(S)^2.$$

The variance of f is the following

$$\text{Var}[f] = \mathbb{E}_{x \in \{-1, 1\}^n} [f^2] - \mathbb{E}_{x \in \{-1, 1\}^n} [f]^2.$$

It is clear that

$$\text{Var}[f] \leq \sum_i^n I_i(f).$$

Below is a straightforward corollary of the above inequality.

Fact 2.3. *If $\text{Var}[f] = 2^{-o(n)}$, then*

$$I_\infty(f) = 2^{-o(n)}.$$

Martingales. Recall that a discrete-time martingale is a sequence of random variables X_0, X_1, X_2, \dots , that satisfies

- For each $n = 0, 1, 2, \dots$, $\mathbb{E}[|X_i|] < \infty$.
- For any $m < n$, $\mathbb{E}[X_m | X_n] = X_n$.

A continuous-time martingale is a stochastic process $(X_t)_{t \geq 0}$ such that

- For any t , $\mathbb{E}[|X_t|] < \infty$.
- For any $s < t$, $\mathbb{E}[X_t | X_s] = X_s$.

A submartingale is a stochastic process with the second property from the above definition replaced by

$$\mathbb{E}[X_t | X_s] \geq X_s.$$

Fact 2.4. *Let X_t, Y_t be martingales.*

- (i) *$aX_t + bY_t$ and $X_t \cdot Y_t$ are also martingales for any constant a, b . Hence any multilinear function of martingales is a martingale.*
- (ii) *If $f : \mathbb{R} \rightarrow \mathbb{R}$ is a convex function, then the process $f(X_t)$ is a submartingale.*

The stopping time τ of a stochastic process is a random variable such that the event $\{\tau \leq t\}$ is completely determined by $X_{\leq t}$. Given two stopping times τ_1, τ_2 , let $\tau_1 \wedge \tau_2$ denote the new stopping time $\min\{\tau_1, \tau_2\}$. For martingales, we have the optional stopping theorem.

Theorem 2.5 (Stopping Theorem). *If τ is almost surely bounded, then*

$$\mathbb{E}[X_\tau] = \mathbb{E}[X_0].$$

For submartingales, the equality is replaced by a greater-than inequality. Finally, the following inequalities will be useful for us.

Theorem 2.6 (Doob’s martingale inequality). *Let X be a submartingale taking real values. Then for any constant $C \geq 0$,*

$$\mathbb{P}\left[\sup_{0 \leq t \leq T} X_t \geq C\right] \leq \frac{\mathbb{E}[\max\{X_T, 0\}]}{C}.$$

Theorem 2.7 (Concentration inequality [5, Theorem 2.21]). *Let X_1, X_2, \dots, X_n be martingales with filtration \mathcal{F} , such that for $i = 1, 2, \dots, n$*

$$\begin{aligned} \text{Var}[X_i \mid \mathcal{F}_{i-1}] &\leq \sigma_i^2, \\ X_i - X_{i-1} &\leq M. \end{aligned}$$

Then

$$\mathbb{P}[X \geq \lambda] \leq \exp\left(-\frac{\lambda^2}{2 \sum \sigma_i^2 + 2M\lambda/3}\right).$$

Finally, we should warn the readers that in this paper, often X is a vector and the subscripts are used for coordinates. In that case, the random process X is denoted $X(t)$, and $X_i(t)$ denotes the evolution of each individual coordinate.

3. CONTROLLED PROCESS

Fix a function $f : \{-1, 1\}^n \rightarrow \{0, 1\}$, and we view $f : \mathbb{R}^n \rightarrow \mathbb{R}$ by considering its multilinear extension. We assume that f is not a constant function. Therefore $f(0) > 0$. In this section, we will consider three different discrete random processes. The first one is the uniform process $X(t) \in \{-1, 0, 1\}^n$ for $t \in \{0, 1, \dots, n\}$. It’s called the uniform process because $X(n)$ will be a uniformly random string from $\{-1, 1\}^n$. The second process $Y(t)$ is obtained from $X(t)$ by conditioning on $f(X(n)) = 1$. Therefore, we call $Y(t)$ the conditioned process. The third process is in effect the same as the second process. They have identical distributions. However, we will take the control theory perspective, and give a player a small number of random coordinates to control. We show that the player will be able to alter a process to the conditioned process, which is otherwise the uniform process. Therefore we sometimes call the third process the controlled process.

First, we consider the following uniform processes $X(t)$ for $t = 0, 1, 2, \dots, n$, such that $X(t) \in \{-1, 0, 1\}^n$, and $X(0) = 0^n$.

Procedure 1 (To generate the discrete uniform process $X(t)$):

Sample a uniformly random permutation $\pi : [n] \rightarrow [n]$.

For time $t = 1, 2, \dots, n$

- Let $i = \pi(t)$. Set $X_i(t)$ to be -1 or 1 uniformly at random.
- For all $j \in [n] \setminus \{i\}$, set $X_j(t) := X_j(t-1)$.

Clearly, the above process is just another way to sample a random element from $\{-1, 1\}^n$. We use the notation P to denote the probability measure over the paths of the above process. The subscript P will be used to emphasize the underlying process and the corresponding measure. For example, $\mathbb{E}_P[f], \mathbb{P}_P[\mathcal{E}]$ are the expectation of the function f and the probability of the event \mathcal{E} , respectively, both defined over the space of the paths of the above process $X(t)$. A crucial component of our analysis is that all the partial derivatives of $f(X(t))$ will be small with high probability for t even very close to n . We formulate it as the following lemma, whose proof requires some technical preparations, and is therefore deferred to Section 5.4.

Lemma 3.1. *Let $\epsilon > 0$ ⁷ be such that*

$$\frac{16}{\epsilon} \ln \frac{4}{\epsilon} \leq \ln \frac{1}{I_\infty(f)}.$$

Then for any $\theta \in (0, 1)$,

$$\mathbb{P}_P \left[\max_{0 \leq t \leq (1-\epsilon)n} |\partial_i f(X(t))| \geq \theta \right] \leq \theta^{-3} I_\infty(f)^{\frac{\epsilon}{16}} + \exp(-\epsilon n/8).$$

Next, we modify Procedure 1 to generate what we call the conditioned process. The goal is to guarantee that the new process ends up being a random element sampled from $f^{-1}(1)$. We use $Y(t)$ to distinguish this new process from $X(t)$. Let Q be a new probability measure defined by the equation

$$\mathbb{P}_Q[Y_i(t) = \pm 1 \mid Y(t-1), \pi(t)] := \frac{1}{2} \pm \frac{\partial_i f(Y(t-1))}{2f(Y(t-1))}. \quad (3.1)$$

A calculation shows that the Radon-Nykodym derivative of the two measures satisfies that for any realization $y(1), y(2), \dots, y(s) \in \{-1, 0, 1\}^n$ of the process $Y(t)$ up to time s ,

$$\begin{aligned} \frac{dQ((y(t))_{1 \leq t \leq s})}{dP((y(t))_{1 \leq t \leq s})} &= \prod_{t=1}^s 2\mathbb{P}_Q[Y_{\pi(t)}(t) = y_{\pi(t)}(t) \mid Y(t-1) = y(t-1)] \\ &= \prod_{t=1}^s \left(1 + y_{\pi(t)}(t) \frac{\partial_t f(y(t-1))}{f(y(t-1))} \right) \\ &= \prod_{t=1}^s \frac{f(y(t))}{f(y(t-1))} \\ &= \frac{f(y(s))}{f(0)}. \end{aligned} \quad (3.2)$$

By taking $s = n$ above, we see that the process $Y(t)$ according to Q is equivalent to the same process $X(t)$ according to P , only conditioned on the event that $f(X(n)) = 1$. In particular, according to Q , $Y(n)$ is just a uniformly random element from $f^{-1}(1)$. Further, if we sample $Y(t)$ and let $\mathcal{R}(t)$ be the restriction induced by $Y(t)$, then $(f|_{\mathcal{R}(t)})^{-1}(1)$ is nonempty for any t as long as f is not the constant 0 function. We record this simple but useful observation that Q is a mild change of measure of P .

⁷Throughout this section, let's assume that ϵn is a positive integer.

Claim 3.2. Let \mathcal{E}_t be some event that depends only on the paths of the random process up to time t , e.g., $X(1), X(2), \dots, X(t)$ according to P or $Y(1), Y(2), \dots, Y(t)$ according to Q . Then for any $t \in [n]$,

$$\mathbb{P}_Q[\mathcal{E}_t] \leq \frac{\mathbb{P}_P[\mathcal{E}_t]}{f(0)}.$$

Proof. This is immediate from (3.2),

$$\mathbb{P}_Q[\mathcal{E}_t] = \mathbb{E}_P \left[\mathbb{I}\{\mathcal{E}_t\} \cdot \frac{dQ}{dP} \right] \leq \frac{\mathbb{P}_P[\mathcal{E}_t]}{f(0)}. \quad \square$$

We summarize the distribution of the “conditioned” process $Y(t)$ according to Q :

Procedure 2 (To generate the conditioned process $Y(t)$):

Sample a uniformly random permutation $\pi : [n] \rightarrow [n]$.

For time $t = 1, 2, \dots, n$:

- Let $i = \pi(t)$. Set $Y_i(t)$ according to the following distribution

$$\mathbb{P}[Y_i(t) = \pm 1] = \frac{1}{2} \pm \frac{\partial_i f(Y(t-1))}{2f(Y(t-1))},$$

- For all $j \in [n] \setminus \{i\}$, set $Y_j(t) := Y_j(t-1)$.

A control-theory point of view. The next step will be to consider the above process as a *controlled version* of the conditioned process $Y(t)$. Fix $\epsilon > 0$ and consider the control problem where at each time t , with probability $1 - \epsilon$, $Y(t)$ does a uniformly random step (according to Procedure 1), and with probability ϵ a player gets to determine the sign of $Y_{\pi(t)}$ according to her own choosing.

The key observation of this section is that as long as the player can control a small fraction of random coordinates, she can simulate the conditioned process exactly. The motivation to study this controlled version of $Y(t)$ is the following: The randomly fixed coordinates out of the player’s control induces a random restriction of the function f . If the player can assign the values to the coordinates of her control, that is the alive coordinates of the corresponding random restriction, to end up in $f^{-1}(1)$, this means the restricted function has a nonempty preimage of 1.

To this end, we consider the following procedure (see [Procedure II](#)) that generates the conditioned process $Y(t)$ as well as the uniform process $X(t)$.

The Procedure II starts with a sampling subroutine as the preparation stage, then followed by two phases that generate $Y(t)$ for time t from 0 to n . The first phase corresponds to that described in the first paragraph of this section. During this phase the player needs to cherish her rare opportunity and play “aggressively.” The second phase starts at a point of time τ when the aggressive strategy no longer works. However, we will show that τ is very close n with high probability. As a result, it would not be a problem to give the player full control and let her play “safely” till the end.

Procedure II (The controlled version of processes $Y(t)$ and $X(t)$):

Sampling Subroutine

- Sample a uniformly random permutation $\pi : [n] \rightarrow [n]$.
- Sample a set $T \subseteq \{1, 2, \dots, n\}$, such that for each $i \in [n]$,

$$\mathbb{P}[i \in T] = \epsilon.$$

T will be the set of times when the player gets to determine the value of the coordinate.

- Sample a uniformly random $z \in \{-1, 1\}^{\pi T}$, the random assignment to the variables not controlled by the player.

Phase 1

Set $Y(0) = X(0) = 0^n$.

For time $t = 1, 2, \dots, n$:

- Let $i = \pi(t)$.
- (Coordinate picked uniformly) If $t \notin T$, set $Y_i(t) = X_i(t) = z_i$.
- (Coordinate determined by player) If $t \in T$, set $Y_i(t)$ and $X_i(t)$ according to the following distributions

$$\mathbb{P}[Y_i(t) = \pm 1] = \frac{1}{2} \pm \frac{1}{2\epsilon} \cdot \frac{\partial_i f(Y(t-1))}{f(Y(t-1))},$$

$$\mathbb{P}[X_i(t) = \pm 1] = \frac{1}{2}.$$

- For all $j \in [n] \setminus \{i\}$, set $Y_j(t) := Y_j(t-1)$, $X_j(t) = X_j(t-1)$.
- If either of the following holds, **exit** this loop

$$\max_{i \in [n]} |\partial_i f(Y(t))| > \epsilon\delta, \quad f(Y(t)) < \delta. \quad \# \text{ the breaking condition}$$

Phase 2

While $t < n$:

- $t = t + 1$.
- Let $i = \pi(t)$. Set $Y_i(t)$ and $X_i(t)$ according to the following distributions

$$\mathbb{P}[Y_i(t) = \pm 1] = \frac{1}{2} \pm \frac{\partial_i f(Y(t-1))}{2f(Y(t-1))},$$

$$\mathbb{P}[X_i(t) = \pm 1] = \frac{1}{2}.$$

- For all $j \in [n] \setminus \{i\}$, set $Y_j(t) := Y_j(t-1)$, $X_j(t) := X_j(t-1)$.

Output $\{Y(t)\}_{t \in \{0, 1, \dots, n\}}, \{X(t)\}_{t \in \{0, 1, \dots, n\}}$.

The process $Y(t)$ will be the main process with which our analysis concerns, whereas the process $X(t)$ is only defined for the sake of entropy comparison: We will later argue that the KL-divergence between the two processes is not too large. It is evident that in Phase 1 the distribution of $Y(1), Y(2), \dots$ according to Procedure II is identical to its distribution according to measure Q as long as

$$|\partial_i f(Y(t-1))| \leq \epsilon f(Y(t-1)). \quad (3.3)$$

Indeed, if (3.3) holds, we have

$$\begin{aligned}\mathbb{P}[Y_{\pi(t)}(t) = \pm 1] &= (1 - \epsilon)\frac{1}{2} + \epsilon \left(\frac{1}{2} \pm \frac{1}{2\epsilon} \cdot \frac{\partial_i f(Y(t-1))}{f(Y(t-1))} \right) \\ &= \frac{1}{2} \pm \frac{\partial_i f(Y(t-1))}{2f(Y(t-1))}.\end{aligned}$$

Let time τ be $n + 1$ if the breaking condition is never hit, otherwise let τ be the time when the breaking condition is hit. The reader may wonder that a more natural choice of the “breaking” condition would be the violation of (3.3). Our definition forces that (i) $f(Y(\tau - 1))$ is large, in addition to that (ii) all derivatives $|\partial_i f(Y(\tau - 1))|$ is small compared to the magnitude of $f(Y(\tau - 1))$. Both facts will be very useful in later sections. Formally, we summarize our definition of τ as below,

$$\tau = \tau_1 \wedge \tau_2 \wedge (n + 1), \quad (3.4)$$

where

$$\begin{aligned}\tau_1 &= \min\{t : \max_{i \in [n]} |\partial_i f(Y(t))| > \epsilon\delta\}. \\ \tau_2 &= \min\{t : f(Y(t)) < \delta\}.\end{aligned}$$

The values of the parameters ϵ, δ will be specified later on. By definition, the condition (3.3) holds for $t < \tau$. We should think τ as a stopping time of Phase 1. After the stopping time τ , the player gets to control each coordinates left. She simply assigns the values according Q as in Procedure 2. Since in both phases Procedure II has the same law as that of Q , the process $Y(t)$ defined in procedure II is identical in distribution to the conditioned process $Y(t)$ defined in Procedure 2. The same is clearly true for the uniform process $X(t)$ in its two versions (Procedure 1 and Procedure II).

In the preparation stage, Procedure II samples a random permutation π , a set T of times controlled by the player and z the random assignment to the variables not controlled by the player. For every $m \in [n]$, let \mathcal{G}_m be the σ -algebra generated by

$$\pi|_{\{1, 2, \dots, m\}}, \quad T \cap \{1, 2, \dots, m\}, \quad \text{and } z|_{\pi\{1, 2, \dots, m\}}.$$

Thus, \mathcal{G}_m contains all the information in a run of Procedure II, excluding the player’s choices, up to time m . Also, \mathcal{G}_m induces a restriction of T :⁸

$$\mathcal{R} = (\pi(\{1, 2, \dots, m\} \setminus T), z).$$

A moment’s thought reveals that run Procedure II, if the controlled process $Y(t)$ satisfies that $\tau > m$, then $f|_{\mathcal{R}}$ contains a nonempty preimage of 1.

Claim 3.3. If $\mathbb{P}[\tau > m \mid \mathcal{G}_m] > 0$, then $(f|_{\mathcal{R}})^{-1}(1) \neq \emptyset$.

Therefore, if we can argue that $\tau > m$ running Procedure II on f and $1 - f$ with the same \mathcal{G}_m , then we actually proved that $f|_{\mathcal{R}}$ is nonconstant. To give a quantitative bound on the variance of $f|_{\mathcal{R}}$ requires some more work. This discussion sets two tasks for the remainder of this section. First, to analyze the stopping time τ and second, to provide the necessary tools to bound the variance of the restricted function.

⁸We can also consider any restriction $\mathcal{R} = (S, z)$ for $S \subseteq \pi(\{1, 2, \dots, m\} \setminus T)$. We will use this observation.

3.1. Stopping time τ of the process $Y(t)$. Next, we prove that with high probability $\tau > (1 - \epsilon)n$ for very small ϵ . Therefore, Phase 2 in Procedure II can not be too long.

Lemma 3.4 (Stopping time τ of the process $Y(t)$). *Let $f : \{-1, 1\}^n \rightarrow \{0, 1\}$ be such that $\text{Var}[f] \geq 2^{-o(n)}$. Further, let $\epsilon > 0$ and δ be such that*

$$\begin{aligned} \frac{16}{\epsilon} \ln \frac{4}{\epsilon} &\leq \ln \frac{1}{I_\infty(f)}, \\ \delta &\geq \frac{I_\infty(f)^{\epsilon/80}}{\epsilon}. \end{aligned}$$

Then for sufficiently large n , we have

$$\mathbb{P}_Q[\tau \leq (1 - \epsilon)n] \leq \frac{3\delta}{f(0)}.$$

Proof. The proof relies on the fact that Q is a mild change of measure with respect to P . Consider the following two bad events,

$$\begin{aligned} \mathcal{E}_1 : \quad \tau_1 &\leq (1 - \epsilon)n, \\ \mathcal{E}_2 : \quad \tau_2 &\leq (1 - \epsilon)n. \end{aligned}$$

We first bound $\mathbb{P}_Q[\mathcal{E}_1]$. Note that

$$\begin{aligned} \mathbb{P}_P[\mathcal{E}_1] &= \mathbb{P}_P \left[\max_{0 \leq s \leq (1-\epsilon)n} |\partial_i f(X(s))| \geq \epsilon\delta \right] \\ &\leq \mathbb{P}_P \left[\max_{0 \leq s \leq (1-\epsilon)n} |\partial_i f(X(s))| \geq I_\infty(f)^{\epsilon/60} \right] \\ &\leq I_\infty(f)^{\epsilon/80} + \exp(-\epsilon n/8), \end{aligned}$$

where the second step holds as $\epsilon\delta \geq I_\infty(f)^{\epsilon/80} \geq I_\infty(f)^{\epsilon/60}$; the final step applies Lemma 3.1. The above bound in turn by Claim 3.2 implies that

$$\begin{aligned} \mathbb{P}_Q[\mathcal{E}_1] &\leq \frac{\mathbb{P}_P[\mathcal{E}_1]}{f(0)} \\ &\leq \frac{I_\infty(f)^{\epsilon/80} + \exp(-\epsilon n/8)}{f(0)}. \end{aligned} \tag{3.5}$$

Next we move to bound $\mathbb{P}_Q[\mathcal{E}_2]$. It is immediate from Claim 3.2:

$$\mathbb{P}_Q[\mathcal{E}_2] \leq \frac{\delta}{f(0)}. \tag{3.6}$$

Apply union bound to (3.5) and (3.6), then for large enough n ,

$$\begin{aligned} \mathbb{P}_Q[\tau \leq (1 - \epsilon)n] &= \mathbb{P}_Q[\mathcal{E}_1 \vee \mathcal{E}_2] \\ &\leq \frac{1}{f(0)} \cdot \left(\delta + I_\infty(f)^{\epsilon/80} + \exp(-\epsilon n/8) \right) \\ &\leq \frac{3\delta}{f(0)} \end{aligned}$$

where in the final step we have $\delta \geq I_\infty(f)^{\epsilon/80} = \exp(-o(\epsilon n))$ by Fact 2.3. \square

3.2. The KL-divergence between $Y(t)$ and $X(t)$. The purpose of this section is to show that for any $m \leq n$, $Y(n)$ given \mathcal{G}_m is close to uniform with high probability over the random choices associated with \mathcal{G}_m . In particular, we will show that the KL-divergence between $Y(n)$ and $X(n)$ given \mathcal{G}_m is small with high probability over the random choices associated with \mathcal{G}_m . Recall that the coordinates of $X(n)$ not fixed by \mathcal{G}_m are uniform.

Lemma 3.5. *For any $m \in [n]$, abbreviate*

$$\mathcal{G}_m = (\pi|_{\{1,2,\dots,m\}}, T \cap \{1,2,\dots,m\}, z|_{\pi\{1,2,\dots,m\}}).$$

Then for some universal constant C , and ϵ, δ in the breaking condition in [Procedure II](#),

$$\mathbb{P}_{\mathcal{G}_m} \left[\text{KL} \left(\frac{Y(n) | \mathcal{G}_m}{X(n) | \mathcal{G}_m} \right) \geq \frac{C}{\epsilon} \ln \frac{en}{n-m+1} \log \frac{e}{\delta} \right] \leq \delta.$$

Proof. Let $\tau' = \tau \wedge (m+1)$. By definition of the stopping time τ (3.4), for $t < \tau'$,

$$f(Y(t)) \geq \delta, \tag{3.7}$$

$$|\partial_i f(Y(t))| \leq \epsilon f(Y(t)). \tag{3.8}$$

We calculate the KL-divergence between $Y(n) | \mathcal{G}_m$ and $X(n) | \mathcal{G}_m$. By the chain rule,

$$\begin{aligned} \text{KL} \left(\frac{Y(n) | \mathcal{G}_m}{X(n) | \mathcal{G}_m} \right) &= \mathbb{E}_{\mathcal{G}_m} \left[\sum_{t=1}^{\tau'-1} \mathbb{I}\{t \in T\} \text{KL} \left(\frac{Y_{\pi(t)}(t) | Y(t-1)}{X_{\pi(t)}} \right) \right. \\ &\quad \left. + \text{KL} \left(\frac{Y(n)|_{\pi\{\tau',\tau'+1,\dots,n\}} | Y(\tau'-1)}{X(n)|_{\pi\{\tau',\tau'+1,\dots,n\}}} \right) \right], \end{aligned}$$

where the equality holds because for any $t \in T$,

$$(Y_{\pi(t)}(t) | Y(t-1)) = (Y_{\pi(t)}(t) | Y(t-1), \mathcal{G}_m),$$

namely, any variable $Y_{\pi(t)}(t)$ controlled by the player is independent of the variables in the future that she has no control of; and all coordinates in $X(n)$ are independent.

Next, using formula (3.2), combined with (3.7), it follows that for any $t \leq \tau'$,

$$\begin{aligned} &\text{KL} \left(\frac{(Y(n)|_{\pi\{t,t+1,\dots,n\}} | Y(t-1))}{X(n)|_{\pi\{t,t+1,\dots,n\}}} \right) \\ &= \log \frac{dQ((Y(i))_{t \leq i \leq n})}{dP((X(i))_{t \leq i \leq n})} = \log \frac{1}{f(Y(t-1))} \leq \log \frac{1}{\delta}. \end{aligned} \tag{3.9}$$

Combining the above two displays,

$$\begin{aligned}
& \text{KL}\left(\frac{Y(n) \mid \mathcal{G}_m}{X(n) \mid \mathcal{G}_m}\right) - \log \frac{1}{\delta} \\
& \leq \sum_{t=1}^n \mathbb{E}_{Y(t-1) \mid \mathcal{G}_m} \left[\mathbb{I}\{t \in T\} \mathbb{I}\{t < \tau'\} \text{KL}\left(\frac{Y_{\pi(t)}(t) \mid Y(t-1)}{X_{\pi(t)}(t)}\right) \right] \\
& = \sum_{t=1}^n \mathbb{E} \left[\mathbb{I}\{t \in T\} \mathbb{I}\{t < \tau'\} \left(1 - H\left(\frac{1}{2} + \frac{\partial_{\pi(t)} f(Y(t-1))}{2\epsilon f(Y(t-1))}\right)\right) \right] \\
& \leq \sum_{t=1}^n \mathbb{E} \left[\mathbb{I}\{t \in T\} \mathbb{I}\{t < \tau'\} \frac{1}{\epsilon^2} \left(\frac{\partial_{\pi(t)} f(Y(t-1))}{f(Y(t-1))}\right)^2 \right],
\end{aligned}$$

where the second step is by the definition of the KL-divergence; and the final step is due to (2.1). Abbreviate

$$Z_t := \mathbb{I}\{t \in T\} \mathbb{I}\{t < \tau'\} \left(\frac{\partial_{\pi(t)} f(Y(t-1))}{f(Y(t-1))}\right)^2.$$

Claim 3.6. There is some universal constant $C \geq 1$, such that for any $t < \tau'$,

$$Z_t \mid Y(t-1) \in [0, \epsilon^2], \quad (3.10)$$

$$\mathbb{E}[Z_t \mid Y(t-1)] \leq \frac{C\epsilon}{n-t+1} \log \frac{e}{\delta}, \quad (3.11)$$

$$\text{Var}[Z_t \mid Y(t-1)] \leq \epsilon^2 \mathbb{E}[Z_t \mid Y(t-1)]. \quad (3.12)$$

Proof. (3.10) follows from (3.8). Let

$$v(t) := \frac{(\nabla f(Y(t-1)))|_{\pi\{t, t+1, \dots, n\}}}{f(Y(t-1))}.$$

Then,

$$\begin{aligned}
\mathbb{E}[Z_t \mid Y(t-1)] &= \mathbb{E}_{\pi(t)} [\epsilon v(t)^2_{\pi(t)} \mid Y(t-1)] \\
&= \frac{\epsilon |v(t)|^2}{n-t+1} \\
&\leq \frac{C\epsilon}{n-t+1} \log \frac{e}{f(Y(t-1))} \\
&\leq \frac{C\epsilon}{n-t+1} \log \frac{e}{\delta},
\end{aligned}$$

where the first step holds as $t \in T$ with probability ϵ ; in the second step, $\pi(t)$ is random within the $n-t+1$ alive coordinates given $Y(t-1)$; the third step follows the Level-1 inequality of Theorem 2.2; and the final step is due to (3.7).

The variance of $Z_t \mid Y(t-1)$ can be bounded as follows:

$$(\epsilon^2 - \mathbb{E}[Z_t]) \mathbb{E}[Z_t] - \text{Var}[Z_t] = \mathbb{E}[(\epsilon^2 - Z_t)Z_t] \geq 0.$$

We comment that such a bound is sometimes referred to as the Bhatia-Davis inequality. \square

By (3.11)-(3.12), the definition that $\tau' \leq m + 1$, and the following elementary fact that

$$\ln(n+1) \leq \sum_{i=1}^n \frac{1}{n} \leq \ln en,$$

we have

$$\sum_{t=1}^n \mathbb{E}[Z_t \mid Y(t-1)] \leq \lambda, \quad (3.13)$$

$$\sum_{t=1}^n \text{Var}[Z_t \mid Y(t-1)] \leq \epsilon^2 \lambda, \quad (3.14)$$

where

$$\lambda = C\epsilon \ln \frac{en}{n-m+1} \log \frac{e}{\delta}.$$

The lemma is concluded by estimating,

$$\begin{aligned} \mathbb{P} \left[\text{KL} \left(\frac{Y(n) \mid \mathcal{G}_m}{X(n) \mid \mathcal{G}_m} \right) \geq \frac{3\lambda}{\epsilon^2} + \log \frac{1}{\delta} \right] \\ \leq \mathbb{P} \left[\sum_{t=1}^n Z_t \mid Y(t-1) \geq 3\lambda \right] \\ \leq \exp \left(-\frac{(2\lambda)^2}{2\epsilon^2 \lambda + 4\epsilon^2 \lambda / 3} \right) \\ \leq \exp \left(-\frac{C}{\epsilon} \ln \frac{en}{n-m+1} \log \frac{e}{\delta} \right) \leq \delta, \end{aligned}$$

where the second step invokes the concentration inequality of Theorem 2.7 since $Z_t - \mathbb{E}[Z_t \mid Y(t-1)]$ is a martingale with respect to $Y(0), Y(1), \dots, Y(t-1)$. This finishes our proof to Lemma 3.5 with a change of the constant C . \square

4. PROOFS OF THE MAIN RESULTS

In this section, we prove the nonasymptotic version of Theorem 1.1, its optimality, and discuss its applications to block sensitivity and decision tree complexity.

Theorem 4.1 (Ain’t over till it’s over). *There are absolute constant $C > 1$. Given $f : \{-1, 1\}^n \rightarrow \{0, 1\}$, such that $I_\infty(f) < 1/C$ and $\text{Var}[f] = 2^{-o(n)}$. Let \mathcal{R} be a random restriction that keeps exactly $\lceil \rho n \rceil$ variables alive, where*

$$\frac{C}{\text{Var}[f]} \cdot \frac{\ln \ln(1/I_\infty(f))}{\ln(1/I_\infty(f))} \leq \rho \leq 1.$$

Let p be such that

$$\frac{8I_\infty(f)^{\rho/C}}{\rho \text{Var}[f]} \leq p \leq 1. \quad (4.1)$$

Then for large enough n ,

$$\mathbb{P} \left[\text{Var}[f|_{\mathcal{R}}] \leq \exp \left(-\frac{C}{\rho} \ln \frac{e}{\rho} \cdot \log \frac{8e}{p \text{Var}[f]} \right) \right] \leq p. \quad (4.2)$$

Before proving the above theorem, we record the following simple fact.

Fact 4.2. *Let $f : \{-1, 1\}^n \rightarrow \{0, 1\}$ be a Boolean function. Let μ be the uniform distribution and γ be some arbitrary distribution over $\{-1, 1\}^n$. If*

$$\gamma(f) \geq \delta, \quad \text{KL}(\gamma \parallel \mu) \leq K.$$

Then

$$\mu(f) \geq 2^{-(K+H(\delta))/\delta}.$$

In particular, if $\gamma(f) = 1$, then

$$\mu(f) \geq 2^{-K}.$$

Proof. Assume that $\gamma(f) = \delta$, and $\text{KL}(\gamma \parallel \mu) = K$. This is without loss of generality because $2^{-(K+H(\delta))/\delta}$ is decreasing in K and increasing in δ by elementary calculus. Let γ_0, γ_1 be the uniform distributions over $f^{-1}(0)$ and $f^{-1}(1)$, respectively. Note $\delta\gamma_1 + (1-\delta)\gamma_0 = \mathbb{E}_\pi[\gamma \circ \pi]$, where π is taken over the product of permutations on $f^{-1}(0)$ and $f^{-1}(1)$. Thus by convexity,

$$\text{KL}(\delta\gamma_1 + (1-\delta)\gamma_0 \parallel \mu) \leq \text{KL}(\gamma \parallel \mu).$$

Consequently, let $\eta = \mu(f)$, then

$$\begin{aligned} & \delta \log \frac{\delta}{\eta} + (1-\delta) \log \frac{1-\delta}{1-\eta} \leq K \\ \implies & \delta \log \frac{1}{\eta} + (1-\delta) \log \frac{1}{1-\eta} \leq K + H(\delta) \\ \implies & \delta \log \frac{1}{\eta} \leq K + H(\delta) \\ \implies & \eta \geq 2^{-(K+H(\delta))/\delta}. \end{aligned}$$

□

Next, we set forth to prove Theorem 4.1. Set

$$\epsilon = \max \left\{ \eta : \eta \leq \frac{\rho}{3}, \eta n \text{ is an integer} \right\}, \quad (4.3)$$

$$\delta = p\text{Var}[f]/8. \quad (4.4)$$

It’s straightforward to verify that for large enough n , and large enough constant C , we have

$$\frac{16}{\epsilon} \ln \frac{4}{\epsilon} \leq \ln \frac{1}{\text{I}_\infty(f)}, \quad (4.5)$$

$$\delta \geq \frac{\text{I}_\infty(f)^{\epsilon/80}}{\epsilon}. \quad (4.6)$$

We will run Procedure II described in Section 3.2 with parameters ϵ and δ . Recall that Procedure II first samples the random permutation π , the set T of times controlled by the player and $z \in \{-1, 1\}^{\pi T}$ is the random assignment for $t \notin T$ in Phase 1. Let $m = (1-\epsilon)n$, $U = \{1, 2, \dots, m\} \setminus T$. Note that by a Chernoff bound, the probability that $|U|$ is less than $\lfloor (1-\rho)n \rfloor$ is at most $\exp(-\epsilon n/2)$. Conditioning on that $|U| \geq \lfloor (1-\rho)n \rfloor$, we randomly sample a set S of $\lfloor (1-\rho)n \rfloor$ elements from U .

Consider the following event

$$\mathcal{E} := \{\tau > m\} \cap \{|U| \geq \lfloor (1-\rho)n \rfloor\}.$$

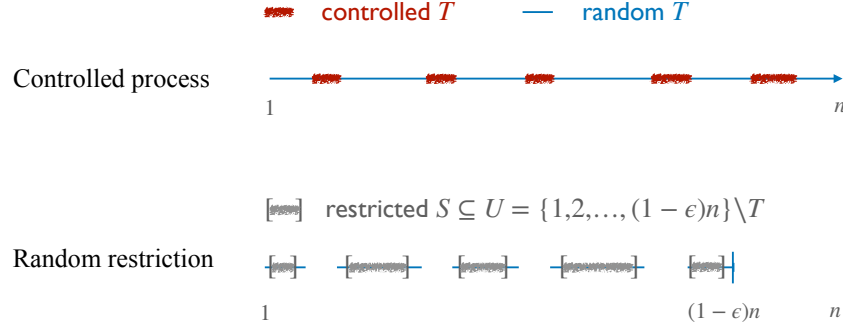


FIGURE 4.1. Random process and random restrictions

The first event in this intersection can be bounded by Lemma 3.4. Hence, by union bound,

$$\begin{aligned} \mathbb{P}[\neg \mathcal{E}] &\leq \frac{3\delta}{f(0)} + \exp(-\epsilon n/2), \\ &\leq \frac{4\delta}{f(0)}, \end{aligned} \quad (4.7)$$

where the second step applies Fact 2.3 with (4.6). Conditioning on \mathcal{E} , $\mathcal{R} = (S, Y(n))$ is distributed as a random restriction that keeps $\lceil \rho n \rceil$ variables alive. Furthermore, the restricted function $f|_{\mathcal{R}}$ satisfies that its mean is bounded away from 0 with high probability.

Claim 4.3. For some universal constant $C' > 1$,

$$\mathbb{P}[\mu(f|_{\mathcal{R}}) < 2^{-K} \mid \mathcal{E}] \leq \frac{2\delta}{\mathbb{P}[\mathcal{E}]}, \quad (4.8)$$

where

$$K = \frac{C'}{\epsilon} \ln \frac{e}{\epsilon} \log \frac{e}{\delta}. \quad (4.9)$$

Our theorem follows immediately from the above claim. Indeed, if \mathcal{E}' is defined analogously to \mathcal{E} where f is replaced by $1 - f$, we have

$$\begin{aligned} \mathbb{P}[\text{Var}[f|_{\mathcal{R}}] < 2^{-K-1}] &\leq \mathbb{P}[\neg \mathcal{E} \vee \neg \mathcal{E}'] + \mathbb{P}[\mu(f|_{\mathcal{R}}) < 2^{-K} \mid \mathcal{E}] + \mathbb{P}[\mu(f|_{\mathcal{R}}) > 1 - 2^{-K} \mid \mathcal{E}'] \\ &\leq \frac{4\delta}{f(0)} + \frac{4\delta}{1 - f(0)} + \frac{2\delta}{\mathbb{P}[\mathcal{E}]} + \frac{2\delta}{\mathbb{P}[\mathcal{E}']} \\ &\leq \frac{8\delta}{\text{Var}[f]} \end{aligned}$$

where in the first step, note that $2^{-K} < 1/2$; the second step plugs in (4.7)-(4.8); in the final step, note that by (4.4),

$$\mathbb{P}[\mathcal{E}], \mathbb{P}[\mathcal{E}'] \geq 1 - \frac{4\delta}{\text{Var}[f]} = 1 - \frac{p}{2} \geq \frac{1}{2} > \text{Var}[f].$$

In view of (4.2) and (4.9), the proof to Theorem 4.1 is finished. It remains to prove Claim 4.3.

Proof of Claim 4.3. Recall that $m = (1-\epsilon)n$ and $U = \{1, 2, \dots, m\} \setminus T$. Abbreviate

$$\mathcal{G}_m = (\pi|_{\{1,2,\dots,m\}}, T \cap \{1, 2, \dots, m\}, z|_{\pi\{1,2,\dots,m\}}),$$

the information of the random process generated by Procedure II excluding the choices of the player up to time m . Further, let γ be the distribution of $Y(n)|_{\pi\bar{U}}$ given \mathcal{G}_m . Then by definition of Procedure II running with respect to function f , $\gamma(f|_{(\pi U, z)}) = 1$. Hence,

$$\begin{aligned} \mathbb{P} [\mu(f|_{(\pi U, z)}) < 2^{-K} \mid \mathcal{E}] \\ &\leq \mathbb{P} \left[\text{KL} \left(\frac{Y(n)}{X(n)} \mid \mathcal{G}_m \right) > K \mid \mathcal{E} \right] \\ &\leq \frac{\delta}{\mathbb{P}[\mathcal{E}]} \end{aligned}$$

where the first step is due to Lemma 4.2; the second step follows Lemma 3.5 for a suitable constant C' in (4.9). Now for any

$$S \in \left(\begin{array}{c} \{1, 2, \dots, m\} \\ \lfloor (1-\rho)n \rfloor \end{array} \right), \quad \text{and } y \in \{-1, 1\}^{\pi S},$$

let $\zeta(S, y)$ be the distribution of $(U, z) \mid \{(S \subseteq U) \wedge (z|_{\pi S} = y)\}$. Then

$$\mu(f|_{(\pi S, y)}) = \mathbb{E}_{(U, z) \sim \zeta(S, y)} [\mu(f|_{(\pi U, z)})].$$

Thus, by Markov's inequality $\mu(f|_{(\pi S, y)}) < 2^{-K-1}$ implies that

$$\mathbb{P}_{(U, z) \sim \zeta(S, y)} [\mu(f|_{(\pi U, z)}) < 2^{-K}] > \frac{1}{2}.$$

Consequently,

$$\begin{aligned} &\frac{1}{2} \mathbb{P}_{\pi, S, y} [\mu(f|_{(\pi S, y)}) < 2^{-K-1}] \\ &\leq \mathbb{P}_{\pi, S, y, (U, z) \sim \zeta(S, y)} [\mu(f|_{(\pi U, z)}) < 2^{-K}] \\ &= \mathbb{P}_{\pi, T, z} [\mu(f|_{(\pi U, z)}) < 2^{-K} \mid \{|U| \geq \lfloor (1-\rho)n \rfloor\}] \\ &\leq \frac{\delta}{\mathbb{P}[\mathcal{E}]} \end{aligned}$$

In view of (4.8), we are done. \square

Remark 4.4. If we consider the random restriction that keeps each variable alive independently with probability ρ , the same statement holds with a slight modification on the proof to the corresponding version of Claim 4.3.

Optimality of our result. Our Theorem 4.1 is essentially optimal with respect to p and ρ . Consider the $(1 - \rho)$ -random restriction $\mathcal{R}_{1-\rho}$. First, we check the optimality in the regime when $\rho = \Omega(1/\log I_\infty(f))$. Consider the majority function $\text{MAJ}_n : \{-1, 1\}^n \rightarrow \{0, 1\}$,

$$\text{MAJ}_n(x) = \begin{cases} 0 & \sum_{i \in n} x_i > 0, \\ 1 & \text{otherwise.} \end{cases}$$

It’s well-known that $I_\infty(\text{MAJ}_n) = \Theta(1/\sqrt{n})$. For $\rho = \Omega(1/\log n)$, let $\mathcal{R}_{1-\rho} = (S, X)$ be the random restriction. Say $|S| = n - k$. With probability at least $1 - \exp(-\Theta(\rho n))$, $k \in (0.5\rho n, 2\rho n)$. Then by Berry-Esseen Theorem, for $\lambda = O(\sqrt{(n - k) \log(n - k)})$,

$$\begin{aligned} \mathbb{P} \left[\sum_{i \in S} X_i \geq \lambda \right] &= \exp \left(-\Theta \left(\frac{\lambda^2}{n - k} \right) \right), \\ \text{Var} \left[\text{MAJ}_n |_{\mathcal{R}_{1-\rho}} \mid \left\{ \sum_{i \in S} X_i \geq \lambda \right\} \right] &\leq \exp \left(-\Theta \left(\frac{\lambda^2}{k} \right) \right). \end{aligned}$$

Thus for $p = \Omega(1/\sqrt{n - k})$,

$$\mathbb{P} \left[\text{Var}[\text{MAJ}_n |_{\mathcal{R}_{1-\rho}}] \leq p^{\Theta(\frac{1}{\rho})} \right] = p.$$

Our bound on the variance is tight up to a $\log(1/\rho)$ factor in the exponent with respect to ρ .

Second, we check the optimality in the regime when $\rho = O(1/\log I_\infty(f))$. Consider the tribes function $\text{TRIBE}_n : \{-1, 1\}^n \rightarrow \{0, 1\}$,

$$\text{TRIBE}_n : x \mapsto \text{AND}_{n/w} \left(\cdots, \bigvee_{j=1}^w x_{ij}, \cdots \right),$$

where for any positive integer w , n is the smallest integral multiple of w such that $\mathbb{P}[\text{TRIBE}_n(x) = 1] \leq 1/2$; $\text{AND}_{n/w} : \{0, 1\}^{n/w} \rightarrow \{0, 1\}$ is the standard logic and function. In particular, $n \approx \ln 2 \cdot w 2^w$, $w = \log n - \log \ln n + o(1)$. Then $I_\infty(\text{TRIBE}_n) = \Theta(\log n/n)$, and $\mu(\text{TRIBE}_n) = \Theta(1)$. Apply random restriction \mathcal{R} that fixes a variable with probability $1 - 1/w = 1 - \Theta(\log 1/I_\infty(\text{TRIBE}_n))$. Then for large enough n ,

$$\mathbb{P}[\text{TRIBE}_n |_{\mathcal{R}} \equiv 1] = \left(1 - \left(\frac{1}{2} + \frac{1}{2w} \right)^w \right)^{n/w} = \Omega(1).$$

Therefore, in this regime $\rho = O(1/\log I_\infty(f))$, with constant probability, there is no variance left under random restrictions for the tribes function. Our bound is tight up to a $\log \log$ factor in the sense that it gives a bound up to the minimum ρ where there is still some variance left after the random restriction.

4.1. Block sensitivity is large almost everywhere. We now move on to our second theorem, concerning block sensitivity. The following is a nonasymptotic version of Theorem 1.2.

Theorem 4.5. *There are absolute constant $C > 1$. For any Boolean function $f : \{-1, 1\}^n \rightarrow \{0, 1\}$, let $\tau = I_\infty(f) < 1/C$, $\text{Var}[f] = 2^{-o(n)}$. Then for large enough n ,*

$$\mathbb{P}_x \left[\text{bs}_f(x) \geq \frac{\text{Var}[f] \ln 1/\tau}{C \ln \ln 1/\tau} \right] \geq 1 - \exp \left(-\Theta \left(\frac{1}{\text{Var}[f]} \ln \ln \frac{1}{\tau} \right) \right).$$

Proof. Let

$$M = \left\lceil \frac{2\text{Var}[f] \ln 1/\tau}{C \ln \ln 1/\tau} \right\rceil.$$

Let $X \in \{-1, 1\}^n$ be random. Randomly partition $[n]$ into M sets, S_1, S_2, \dots, S_M , each of size $\lfloor n/M \rfloor$ with maybe a small number of remaining indices. Note that for any $i \in [M]$, $\mathcal{R} = (S_i, X)$ is a random restriction of fixed size. Then by Theorem 4.1, with probability at least $1 - \exp(-\Theta(\log \log(1/\tau)/\text{Var}[f]))$, $\text{Var}[f|_{\mathcal{R}}] > 0$. In that case, exists $T_i \subseteq S_i$, such that $f(X \oplus (-1)^{\mathbf{1}_{T_i}}) \neq f(X)$. The statement thus holds by the following double-counting principle,

$$\frac{1}{2} \mathbb{P}_x[\text{bs}_f(x) < M/2] \leq \mathbb{P}_{\mathcal{R}}[f|_{\mathcal{R}} \text{ is constant}]. \quad \square$$

4.2. Decision tree complexity of random restriction to monotone functions. We record another application of our main result regarding the decision tree complexity of the restricted function, which is in some sense a reverse statement to the famous Håstad’s switching lemma. Let $\text{DT}(f)$ denote the deterministic decision tree complexity of f .⁹

Theorem 4.6 (decision tree complexity of random restriction). *There are absolute constant $C > 1$. For any monotone function $f : \{-1, 1\}^n \rightarrow \{0, 1\}$, such that*

$$\log \left(\frac{1}{I_\infty(f)} \right) \geq C \log \left(\frac{1}{\text{Var}[f]} \right). \quad (4.10)$$

Let \mathcal{R} be a random restriction that keeps exactly $\lceil \rho n \rceil$ variables alive, where

$$\rho = \Omega \left(\sqrt{\frac{\log \text{Var}[f]}{\log I_\infty(f)}} \log \frac{\log I_\infty(f)}{\log \text{Var}[f]} \right).$$

Then for large enough n ,

$$\mathbb{P} \left[\text{DT}(f|_{\mathcal{R}}) \geq I_\infty(f)^{-\Theta(\rho)} \right] \geq \frac{1}{2}. \quad (4.11)$$

Proof. For simplicity, assume that ρn is a positive integer. We need the following well-known result due to O’Donnell et al. [18]: For any Boolean function h ,

$$I_\infty(h) \cdot \text{DT}(h) \geq \text{Var}[h]. \quad (4.12)$$

Consider the uniform process $X(t)$. By definition, $X((1 - \rho)n)$ induces a random restriction \mathcal{R} that keeps exactly ρn variables alive. By Lemma 3.1,

$$\mathbb{P}_P \left[\max_{0 \leq t \leq (1-\rho)n} |\partial_i f(X(t))| \geq I_\infty(f)^{\frac{\rho}{40}} \right] \leq I_\infty(f)^{\frac{\rho}{40}} + \exp(-\rho n/8).$$

⁹Although the theorem is stated with respect to the deterministic decision tree complexity, one can replace the deterministic decision tree complexity by many other complexity measures, for example, the randomized decision tree complexity, as they are polynomially related for total functions.

Since f is monotone, the influence $I_i(f|\mathcal{R}) = |\partial_i f(X((1-\rho)n))|$ for any alive coordinate i . The above formula then implies

$$\mathbb{P}_{\mathcal{R}} \left[I_{\infty}(f|\mathcal{R}) \geq I_{\infty}(f)^{\frac{\rho}{30}} \right] \leq I_{\infty}(f)^{\frac{\rho}{40}} + \exp(-\rho n/8). \quad (4.13)$$

By Theorem 4.1,

$$\mathbb{P} \left[\text{Var}[f|\mathcal{R}] \leq \exp \left(-\frac{C}{\rho} \log \frac{e}{\rho} \cdot \log \frac{8e}{\text{Var}[f]p} \right) \right] \leq p. \quad (4.14)$$

Set $p = 1/3$. Then combining (4.12)-(4.14),

$$\begin{aligned} & \mathbb{P} \left[\text{DT}(f) \geq I_{\infty}(f)^{-\frac{\rho}{60}} \right] \\ & \geq \mathbb{P} \left[\text{DT}(f) \geq \exp \left(-\frac{C}{\rho} \log \frac{e}{\rho} \cdot \log \frac{8e}{\text{Var}[f]p} + \frac{\rho}{30} \ln \frac{1}{I_{\infty}(f)} \right) \right] \\ & \geq \mathbb{P} \left[\text{Var}[f|\mathcal{R}] \geq \exp \left(-\frac{C}{\rho} \log \frac{e}{\rho} \cdot \log \frac{8e}{\text{Var}[f]p} \right) \wedge \left(I_{\infty}(f|\mathcal{R}) \leq I_{\infty}(f)^{\frac{\rho}{30}} \right) \right] \\ & \geq 1 - p - I_{\infty}(f)^{\frac{\rho}{40}} - \exp(-\rho n/8) \\ & \geq 1 - p - 2I_{\infty}(f)^{\frac{\rho}{40}}, \end{aligned}$$

where the first step holds for our choice of ρ and (4.10); the final step holds by Fact 2.3. Finally, by our choice of ρ and the bound on $I_{\infty}(f)$, we have $2I_{\infty}(f)^{\frac{\rho}{40}} < 1/6$. In view of (4.11), we have finished the proof. \square

5. RANDOM RESTRICTIONS AND HYPERCONTRACTIVITY

In this section, we consider the continuous random process revealing information about the inputs $X \in \{-1, 1\}^n$ gradually in a bit by bit manner. We establish a hypercontractivity theorem for this “operator,” and then use the new hypercontractivity theorem to show that the first-order Fourier coefficients remain small under random restriction given that the original function has small individual influences.

5.1. A martingale setup for random restrictions. Consider the following random process. Let $x \in \{-1, 1\}^n$ be a uniformly random element. Let $(\tau_i)_{i \in [n]}$ be random variables uniformly distributed in the interval $[0, 1]$. τ induces a permutation on $[n]$. This is essentially the only relevant information. For technical reasons we prefer this continuous description in this section. Define $S(t) = \{i : \tau_i \leq t\}$, and define process $X(t) \in [-1, 1]^n$ as follows

$$X_i(t) = \begin{cases} 0 & \tau_i > t, \\ x_i & \tau_i \leq t. \end{cases}$$

In another word, a random ± 1 variable is revealed with probability t at time t . This random process induces a random restriction $\mathcal{R}(t) = (S(t), Y)$ of function f , that all the variables in $S(t)$ is set according to Y while the other variables are left alive. Below, we collect some properties of a function f with respect to the above process.

Proposition 5.1. *For any multilinear function $f : [-1, 1]^n \rightarrow \mathbb{R}$ and any $t \geq 0$,*

- (i) $\mathbb{E}[|\nabla f(X(t))|^2] \leq \|f\|_{\infty}^2/(1-t),$
- (ii) $\mathbb{E}[\partial_i f(X(t))^2] \leq I_i[f], \text{ for } i = 1, 2, \dots, n.$

Proof. (i) Note that for $i \notin S(t)$, by definition

$$\partial_i f(X(t)) = \widehat{f|_{\mathcal{R}(t)}}(i).$$

Thus, by Parseval’s identity,

$$\sum_{i=1}^n \mathbb{I}\{\tau_i > t\} \partial_i f(X(t))^2 \leq \mathbb{E}[(f|_{\mathcal{R}(t)})^2] \leq \|f\|_\infty^2.$$

Since $\mathbb{I}\{\tau_i > t\}$ and $\partial_i f(X(t))^2$ are independent, we have

$$\|f\|_\infty^2 \geq \mathbb{E} \left[\sum_{i=1}^n \mathbb{I}\{\tau_i > t\} \partial_i f(X(t))^2 \right] = (1-t) \mathbb{E}[|\nabla f(X(t))|^2].$$

(ii) By Fourier expansion of $\partial_i f$,

$$\begin{aligned} \mathbb{E} [\partial_i f(X(t))^2] &= \mathbb{E} \left[\left(\sum_{S \ni i} \hat{f}(S) \chi_{S \setminus \{i\}}(X(t)) \right)^2 \right] \\ &= \mathbb{E} \left[\sum_{S \ni i} \hat{f}(S)^2 \mathbb{I}\{\tau_j \leq t, \forall j \in S \setminus \{i\}\} \right] \\ &\leq I_i(f). \end{aligned} \quad \square$$

5.2. A hypercontractive inequality for random restrictions. As the time t increases, the process $X(t)$ reveals more information about the location of $X(1)$. Thus, for $0 \leq t \leq T \leq 1$, we may view $f(X(t))$ as a “noisy” version of $f(X(T))$. It is therefore expected that some hypercontractive inequality holds for those two expressions. This intuition can be made concrete by the following theorem.

Theorem 5.2 (A hypercontractive inequality). *For any $0 \leq t \leq T \leq 1$, and any multilinear $f : [-1, 1]^n \rightarrow \mathbb{R}$, we have for the random process X defined in the previous section,*

$$(\mathbb{E} |f(X(t))|^{2+\epsilon})^{\frac{1}{2+\epsilon}} \leq (\mathbb{E} |f(X(T))|^2)^{1/2}, \quad (5.1)$$

where

$$\epsilon = T - t.$$

Proof. The proof is by induction on n . Once we establish the base case, the inductive step follows from a standard argument. We first show the inductive step since the base case is more involved.

Inductive step. Let $f(z, x) = zg(x) + h(x)$, where $x \in [-1, 1]^n$, and $z \in [-1, 1]$. Then

$$\begin{aligned}
& \mathbb{E}[|f(Z(t), X(t))|^{2+\epsilon}]^{\frac{1}{2+\epsilon}} \\
&= \left(\mathbb{E}_{X(t)} \left[\mathbb{E}_{Z(t)} [|Z(t)g(X(t)) + h(X(t))|^{2+\epsilon}] \right] \right)^{\frac{1}{2+\epsilon}} \\
&\stackrel{(i)}{\leq} \left(\mathbb{E}_{X(t)} \left[\mathbb{E}_{Z(T)} [|Z(T)g(X(t)) + h(X(t))|^2]^{\frac{2+\epsilon}{2}} \right] \right)^{\frac{1}{2+\epsilon}} \\
&\stackrel{(ii)}{\leq} \left(\mathbb{E}_{Z(T)} \left[\mathbb{E}_{X(t)} [|Z(T)g(X(t)) + h(X(t))|^{2+\epsilon}]^{\frac{2}{2+\epsilon}} \right] \right)^{\frac{1}{2}} \\
&\stackrel{(iii)}{\leq} \left(\mathbb{E}_{Z(T)} \left[\mathbb{E}_{X(T)} [(Z(T)g(X(T)) + h(X(T)))^2] \right] \right)^{\frac{1}{2}},
\end{aligned}$$

where (i) holds because for any fixed $X(t)$, $f = z \cdot g(X(t)) + h(X(t))$ is a multilinear function on z , thus we can apply the inductive hypothesis; inequality (iii) is true, again because for any fixed $Z(T)$, f is a multilinear function on x and we apply the inductive hypothesis; (ii) follows by the Minkowski inequality, in particular,

$$\left(\mathbb{E}_x \left[\mathbb{E}_z [f(z, x)^2]^{\frac{2+\epsilon}{2}} \right] \right)^{\frac{2}{2+\epsilon}} \leq \mathbb{E}_z \left[\mathbb{E}_x [|f(z, x)|^{2+\epsilon}]^{\frac{2}{2+\epsilon}} \right].$$

Base case. For the base case we consider two scenarios separately. (i) f is nonnegative (or, nonpositive) function. Let $f : [-1, 1] \rightarrow [0, \infty)$, say $f = ax + b$. It suffices to consider the special case when $f = ax + 1$ for some $0 < a < 1$ after normalization. The reason is as follows: Since f is nonnegative, $0 \leq |a| \leq b$. Thus, we can assume $a \geq 0$, this assumption does not change $\mathbb{E}[|f(X(t))|^p]$. For $b = 0$, there is nothing to prove. So we can assume $b = 1$ by normalization. For $a = 0$, f is constant function. The statement is clearly true. Finally, for the case $a = 1$, it follows from continuity. After the above simplification, we make the actual analysis.

$$\begin{aligned}
\mathbb{E}[(aX(t) + 1)^{2+\epsilon}] &= (1-t) + \frac{t}{2}((1+a)^{2+\epsilon} + (1-a)^{2+\epsilon}) \\
&= 1 - t + t \sum_{k \geq 0} a^{2k} \binom{2+\epsilon}{2k}, \\
&= 1 + t \sum_{k > 0} a^{2k} \binom{2+\epsilon}{2k},
\end{aligned}$$

where the second step uses Taylor expansion of $(1+x)^p$ for $|x| < 1$. Note that for $\epsilon \in [0, 1]$ and any $k \geq 2$,

$$\binom{2+\epsilon}{2k} \leq 0.$$

Hence,

$$\mathbb{E}[(aX(t) + 1)^{2+\epsilon}] \leq 1 + t(1 + \epsilon/2)(1 + \epsilon)a^2. \tag{5.2}$$

On the other hand,

$$\begin{aligned}
\mathbb{E}[(aX(T) + 1)^2]^{\frac{2+\epsilon}{2}} &= (1 + Ta^2)^{\frac{2+\epsilon}{2}} \\
&\geq 1 + (1 + \epsilon/2)Ta^2,
\end{aligned} \tag{5.3}$$

where the last step follows from Fact 2.1. Compare (5.2) and (5.3), we get that

$$\mathbb{E}[(aX(t) + 1)^{2+\epsilon}] \leq \mathbb{E}[(aX(T) + 1)^2]^{1+\epsilon/2}$$

as long as

$$\epsilon \leq \frac{T-t}{t} \quad (5.4)$$

and $\epsilon \in [0, 1]$.

(ii) f takes both positive and negative values, say $f = ax + b$. This time it suffices to consider the special case when $f = x + b$, for $0 < b < 1$. Since if a is not 1, we can consider the function f/a . In addition, changing b to $|b|$ does not affect $\mathbb{E}[|f|^p]$. Then

$$\begin{aligned} \mathbb{E}[|X(t) + b|^{2+\epsilon}] &= (1-t)b^{2+\epsilon} + t/2((1+b)^{2+\epsilon} + (1-b)^{2+\epsilon}) \\ &= (1-t)b^{2+\epsilon} + t \sum_{k \geq 0} b^{2k} \binom{2+\epsilon}{2k} \\ &\leq (1-t)b^2 + t(1+b^2(1+\epsilon/2)(1+\epsilon)) \\ &= 1 + (1-t)(b^2 - 1) + t(1+\epsilon/2)(1+\epsilon)b^2, \end{aligned} \quad (5.5)$$

where the third step uses the facts that $\binom{2+\epsilon}{2k} \leq 0$ for $k \geq 2$ and $\epsilon \in [0, 1]$, and that b^x is decreasing on x for $0 < b < 1$. On the other hand,

$$\begin{aligned} &(\mathbb{E}[(X(T) + b)^2])^{\frac{2+\epsilon}{2}} \\ &= (T + b^2)^{1+\epsilon/2} \\ &= (1 + b^2)^{1+\epsilon/2} \left(1 - \frac{1-T}{1+b^2}\right)^{1+\epsilon/2} \\ &\geq (1 + (1+\epsilon/2)b^2) \left(1 - (1+\epsilon/2)\frac{1-T}{1+b^2}\right) \\ &= 1 + (1+\epsilon/2)b^2 - (1+\epsilon/2)(1+b^2+\epsilon b^2/2)\frac{1-T}{1+b^2} \\ &= 1 + (1+\epsilon/2)b^2 - (1+\epsilon/2)(1-T) - (1+\epsilon/2)b^2\epsilon(1-T)/(2+2b^2) \\ &\geq 1 + (1+\epsilon/2)b^2 - (1+\epsilon/2)(1-T) - (1+\epsilon/2)b^2\epsilon(1-T)/2, \end{aligned} \quad (5.6)$$

where the third step invokes Fact 2.1 twice. Let R, L denote (5.6) and (5.5), respectively. Further, let $B = b^2$, then $R - L$ is a linear function in B . To verify that $R \geq L$, one only needs verify the cases when $B = 0$ and $B = 1$. Recall that $\epsilon = T - t$, therefore

$$\begin{aligned} B = 0: \quad R - L &= 1 - (1+\epsilon/2)(1-T) - t \\ &= \epsilon - (1-T)\epsilon/2 \\ &= \epsilon(1+T)/2 \\ &\geq 0, \end{aligned}$$

and

$$\begin{aligned}
 B = 1 : \quad R - L &= (1 + \epsilon/2)(T - (1 - T)\epsilon/2 - t(1 + \epsilon)) \\
 &= (1 + \epsilon/2)(\epsilon/2 + T\epsilon/2 - \epsilon t) \\
 &= (1 + \epsilon/2)\epsilon/2(1 + T - 2t) \\
 &\geq 0.
 \end{aligned}$$

This concludes our proof. \square

Remark 5.3. One can also prove a hypercontractive inequality of the p -norm vs. 2-norm for $1 < p < 2$. The proof is analogous.

5.3. ℓ_∞ -Fourier mass of $f|_{\mathcal{R}(t)}$ of the first order. A key quantity in our analysis is the ℓ_∞ -Fourier mass of $f|_{\mathcal{R}(t)}$ of the first order. Namely,

$$\beta^*(t) = \max_{i \notin S(t)} |\partial_i f(X(t))|. \quad (5.7)$$

In some sense, $\beta^*(t)$ represents the maximal influence of $f|_{\mathcal{R}(t)}$. In particular, for the special case when f is a monotone Boolean function, $\beta^*(t)$ is exactly $I_\infty(f|_{\mathcal{R}(t)})$. The importance of $\beta^*(t)$ will become clear in later sections. Next, we show that with high probability $\beta^*(t)$ remains small for t even very close to 1. In fact, what we will show is that

$$\beta = \max_{i \in [n]} |\partial_i f(X(t))|$$

remains small with high probability. In particular, we establish the following lemma using the hypercontractive inequality from the last section.

Lemma 5.4 (“influence” remains small under random restriction). *Given $f : \{-1, 1\}^n \rightarrow [-1, 1]$. For any $0 \leq t < 1$ such that*

$$\frac{8}{1-t} \ln \frac{2}{1-t} \leq \ln \frac{1}{I_\infty(f)}. \quad (5.8)$$

Then for any $\theta \in (0, 1)$,

$$\mathbb{P} \left[\sup_{0 \leq s \leq t} \beta(s) \geq \theta \right] \leq \theta^{-3} I_\infty(f)^{\frac{1-t}{8}}.$$

Proof. Take $T = (1 + t)/2$ and let

$$\epsilon = T - t. \quad (5.9)$$

Then

$$\begin{aligned}
\mathbb{P} \left[\sup_{0 \leq s \leq t} \beta(s) \geq \theta \right] &\leq \mathbb{P} \left[\sup_{0 \leq s \leq t} \sum_{i=1}^n |\partial_i f(X(s))|^{2+\epsilon} \geq \theta^{2+\epsilon} \right] \\
&\stackrel{(i)}{\leq} \theta^{-2-\epsilon} \sum_i \mathbb{E}[|\partial_i f(X(t))|^{2+\epsilon}] \\
&\stackrel{(ii)}{\leq} \theta^{-2-\epsilon} \sum_i (\mathbb{E}[\partial_i f(X(T))^2])^{1+\epsilon/2} \\
&\stackrel{(iii)}{\leq} \theta^{-2-\epsilon} \sum_i \mathbb{I}_i(f)^{\epsilon/2} \mathbb{E}[\partial_i f(X(T))^2] \\
&\stackrel{(iv)}{\leq} \theta^{-2-\epsilon} \frac{\mathbb{I}_\infty(f)^{\epsilon/2}}{1-T} \\
&\stackrel{(v)}{\leq} \theta^{-2-\epsilon} \mathbb{I}_\infty(f)^{\epsilon/4}
\end{aligned} \tag{5.10}$$

where (i) is true due to Fact 2.4 and Theorem 2.6; (ii) follows from Theorem 5.2, (iii) follows from Proposition 5.1 (ii), (iv) follows from Proposition 5.1 (i) and (v) follows by our choice of T , and (5.8). \square

5.4. Proof of Lemma 3.1. We have almost proved Lemma 3.1 in the previous section except that we proved the version with the continuous random process instead of the discrete one. Next, we show that the continuous random process and the corresponding probability measure \tilde{P} used in Lemma 5.4 is close to the discrete uniform process generated by Procedure 1 in Section 3 with measure P in the following sense.

Claim 5.5. Let \mathcal{E}_t be some event that depends only on $X(t)$. Then for any $\epsilon \in (0, 1)$,

$$\mathbb{P}_P \left[\bigvee_{0 \leq t \leq (1-\epsilon)n} \mathcal{E}_t \right] \leq \mathbb{P}_{\tilde{P}} \left[\bigvee_{0 \leq \tilde{t} \leq (1-\epsilon/2)} \mathcal{E}_{\tilde{t}} \right] + \exp(-\epsilon n/8).$$

Proof. We couple the two processes in the obvious way. Recall that $(\tau_i)_{i \in [n]}$ is the random variables uniformly distributed in the interval $[0, 1]$ in the continuous process. τ induces a permutation π on $[n]$. As time \tilde{t} goes from 0 to 1 in \tilde{P} , whenever a variable is set to value $v \in \{-1, 1\}$, the corresponding variable in the discrete process is also set to v . Recall that we denote the set of fixed variables at time \tilde{t} in \tilde{P} by $S(\tilde{t})$. Then at time $\tilde{t} = (1 - \epsilon/2)$, by Chernoff bound,

$$\mathbb{P}_{\tilde{P}}[|S(\tilde{t})| < (1 - \epsilon)n] \leq \exp(-\epsilon n/8).$$

Conditioning on $|S(\tilde{t})| \geq (1 - \epsilon)n$,

$$\left\{ \bigvee_{0 \leq t \leq (1-\epsilon)n} \mathcal{E}_t \right\}_P \iff \left\{ \bigvee_{0 \leq \tilde{t} \leq (1-\epsilon/2)} \mathcal{E}_{\tilde{t}} \right\}_{\tilde{P}}.$$

The claim follows. \square

Now, Lemma 3.1 is an immediate corollary of Lemma 5.4 and Claim 5.5.

Corollary 5.6 (Restatement of Lemma 3.1). *Let $\epsilon > 0$ be such that*

$$\frac{16}{\epsilon} \ln \frac{4}{\epsilon} \leq \ln \frac{1}{I_\infty(f)}.$$

Then for any $\theta \in (0, 1)$,

$$\mathbb{P}_P \left[\max_{0 \leq t \leq (1-\epsilon)n} |\partial_i f(X(t))| \geq \theta \right] \leq \theta^{-3} I_\infty(f)^{\frac{\epsilon}{16}} + \exp(-\epsilon n/8).$$

Proof. Let $t = (1 - \epsilon/2)$. The choice of ϵ guarantees that we can apply Lemma 5.4. In view of Claim 5.5, we are done. \square

REFERENCES

- [1] Scott Aaronson and Andris Ambainis. The need for structure in quantum speedups. *Theory of Computing*, 10(6):133–166, 2014. [1](#)
- [2] Ryan Alweiss, Shachar Lovett, Kewen Wu, and Jiapeng Zhang. Improved bounds for the sunflower lemma. *Annals of Mathematics*, 194(3):795 – 815, 2021. [1](#)
- [3] M. Bellare, O. Goldreich, and M. Sudan. Free bits, PCPs and non-approximability-towards tight results. In *Proceedings of IEEE 36th Annual Foundations of Computer Science*, pages 422–431, 1995. [1](#)
- [4] Michael Ben-Or and Nathan Linial. Collective coin flipping, robust voting schemes and minima of banzhaf values. In *26th Annual Symposium on Foundations of Computer Science (sfcs 1985)*, pages 408–416, 1985. [1](#)
- [5] Fan Chung and Linyuan Lu. *Complex graphs and networks*. Number 107 in Conference Board of the mathematical science. American Mathematical Soc., 2006. [2.7](#)
- [6] Thomas M Cover and Joy A Thomas. Elements of information theory 2nd edition (wiley series in telecommunications and signal processing). *Acessado em*, 2006. [5](#)
- [7] Ronen Eldan and Renan Gross. Decomposition of mean-field Gibbs distributions into product measures. *Electronic Journal of Probability*, 23(none):1 – 24, 2018. [1](#)
- [8] E. Friedgut. Boolean functions with low average sensitivity depend on few coordinates. *Combinatorica*, 1(18):27–35, 1998. [1](#)
- [9] Johan Håstad. *Computational limitations of small-depth circuits*. MIT Press, Cambridge, MA, USA, 1987. [1](#)
- [10] Johan Håstad. Testing of the long code and hardness for clique. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing*, pages 11–19, New York, NY, USA, 1996. Association for Computing Machinery. [1](#)
- [11] Johan Håstad. Some optimal inapproximability results. *J. ACM*, 48(4):798–859, jul 2001. [1](#)
- [12] J. Kahn, G. Kalai, and N. Linial. The influence of variables on boolean functions. In *Proceedings of the 29th Annual Symposium on Foundations of Computer Science*, pages 68–80, 1988. [1](#)
- [13] Gil Kalai. A fourier-theoretic perspective on the Condorcet paradox and Arrow’s theorem. *Advances in Applied Mathematics*, 29(3):412–426, 2002. [1](#), [1](#)
- [14] Gil Kalai. Social indeterminacy. *Econometrica*, 72(5):1565–1581, 2004. [1](#), [1](#)
- [15] Subhash Khot, Guy Kindler, Elchanan Mossel, and Ryan O’Donnell. Optimal inapproximability results for MAX-CUT and other 2-variable CSPs? *SIAM Journal on Computing*, 37(1):319–357, 2007. [1](#)
- [16] David Lichtenstein, Nathan Linial, and Michael E. Saks. Some extremal problems arising from discrete control processes. *Comb.*, 9(3):269–287, 1989. [4](#)
- [17] Elchanan Mossel, Ryan O’Donnell, and Krzysztof Oleszkiewicz. Noise stability of functions with low influences: Invariance and optimality. *Annals of Mathematics*, pages 295–341, 2010. [1](#), [1](#)
- [18] R. O’Donnell, M. Saks, O. Schramm, and R.A. Servedio. Every decision tree has an influential variable. In *46th Annual IEEE Symposium on Foundations of Computer Science (FOCS’05)*, pages 31–39, 2005. [1](#), [4.2](#)
- [19] Ryan O’Donnell. Hardness amplification within NP. *Journal of Computer and System Sciences*, 69(1):68–94, 2004. Special Issue on Computational Complexity 2002. [1](#)

- [20] Michel Talagrand. On Russo’s approximate zero-one law. *Ann. Probab.*, 22(3):1576–1587, 1994. [1](#)