

A. Near-Matching Lower Bound: Theorem 3

Our function f witnessing the lower bound of [Theorem 3](#) will be a variant of the well-known $\text{TRIBES}_\ell : \{0, 1\}^\ell \rightarrow \{0, 1\}$ function from the analysis of boolean functions. Recall that TRIBES_ℓ is a read-once DNF formula with $m := \lfloor \frac{\ell}{w} \rfloor$ terms of width exactly w over disjoint variables (with some variables possibly left unused),

$$\text{TRIBES}_\ell(x) := T_1(x) \vee \cdots \vee T_m(x),$$

and $w := \log \ell - \log \ln \ell + o_\ell(1)$ is chosen so that $\Pr[\text{TRIBES}_\ell(\mathbf{x}) = 1]$ is as close to $\frac{1}{2}$ as possible.⁴ (For more on the TRIBES function, see Chapter §4.2 of ([O'Donnell, 2014](#))). Our construction will also involve the Majority function, $\text{MAJ}_k(y) := \mathbb{1}[\sum_{i=1}^k y_i \geq 0]$.

The function that witnesses the separation. Let $m' < m$ be chosen so that the function $\text{TRIBES}'_\ell : \{0, 1\}^\ell \rightarrow \{0, 1\}$,

$$\text{TRIBES}'_\ell(x) := T_1(x) \vee \cdots \vee T_{m'}(x),$$

has acceptance probability $\Pr[\text{TRIBES}'_\ell(\mathbf{x}) = 1]$ as close to 0.499 as possible.⁵ We additionally define $\text{REST}_\ell : \{0, 1\}^\ell \rightarrow \{0, 1\}$ to be:

$$\text{REST}_\ell(x) := \text{TRIBES}_\ell(x) \wedge \neg \text{TRIBES}'_\ell(x),$$

nothing that $\Pr[\text{REST}_\ell(\mathbf{x}) = 1] = 0.001$.

Our function $f_{\ell,k} : \{0, 1\}^\ell \times \{0, 1\}^k \rightarrow \{0, 1\}$ is defined as follows:

$$f_{\ell,k}(x, y) := \begin{cases} 1 & \text{if } \text{TRIBES}'_\ell(x) = 1 \\ \text{MAJ}_k(y) & \text{if } \text{REST}_\ell(x) = 1 \\ 0 & \text{otherwise,} \end{cases}$$

where k and ℓ are chosen so that [Equation \(2\)](#) in the statement of [Lemma A.2](#) below is satisfied with equality.

Proposition A.1 (Upper bound on opt). $\text{opt}_{f_{\ell,k}, 2^\ell} \leq 0.01$.

Proof. We have that

$$\begin{aligned} \Pr[f_{\ell,k}(\mathbf{x}, \mathbf{y}) \neq \text{TRIBES}_\ell(\mathbf{x})] &= \Pr[\text{REST}_\ell(\mathbf{x}) = 1 \text{ and } \text{MAJ}_k(\mathbf{y}) \neq \text{TRIBES}_\ell(\mathbf{x})] \\ &= \Pr[\text{REST}_\ell(\mathbf{x}) = 1] \cdot \Pr[\text{MAJ}_k(\mathbf{y}) = 0] \\ &= 0.001 \cdot \frac{1}{2} < 0.01. \end{aligned}$$

Since TRIBES_ℓ is computed by a decision tree of size $\leq 2^\ell$, it follows that that $\text{opt}_{f_{\ell,k}, 2^\ell} \leq 0.01$. □

Lemma A.2 (Lower bound on TOPDOWNERROR). *There are universal constants c_1 and c_2 such that the following holds. Suppose k and ℓ satisfy:*

$$\frac{c_1}{\sqrt{k}} \geq \frac{\log \ell}{\ell}. \quad (2)$$

Then $\text{TOPDOWNERROR}(f_{\ell,k}, 2^{c_2 k / \log k}) \geq 0.49$.

Proof of Theorem 3 assuming Lemma A.2. We choose k and ℓ such that [Equation \(2\)](#) is satisfied with equality, and define $s := 2^\ell$. [Theorem 3](#) follows as an immediate consequence of [Proposition A.1](#) and [Lemma A.2](#) since $2^{\Omega(k / \log k)} = 2^{\Omega(\ell^2 / (\log \ell)^3)} = s^{\Omega(\log s)}$. □

⁴While this acceptance probability cannot be made exactly $\frac{1}{2}$ for all values of ℓ due to granularity issues, it will be the case that $\Pr[\text{TRIBES}_\ell(\mathbf{x})] = \frac{1}{2} \pm O(\frac{\log \ell}{\ell})$. For clarity we will assume for the remainder of this proof the acceptance probability of TRIBES_ℓ is exactly $\frac{1}{2}$, noting that the same calculations go through if one carries around the additive $o_\ell(1)$ factor.

⁵The same remark as in the previous footnote applies here.

A.1. Proof of Lemma A.2

Our proof of Lemma A.2 draws on many of the ideas in (Blanc et al., 2020)'s proof of their Theorem 6(b). Let T denote the tree constructed by $\text{BUILDTOPDOWNDT}(f_{\ell,k}, 2^{c_2 k / \log k})$, where c_2 is a universal constant that will be determined later.

Claim A.3. *Let $\bar{\pi}$ be a path in T that leads to a first query to an x -variable. Then $\text{Inf}_{y_j}(\text{MAJ}_k(y)_{\bar{\pi}}) \leq \frac{1}{100\sqrt{k}}$ for all $j \in [k]$.*

Proof. Suppose without loss of generality that $\bar{\pi}$ leads to a query to x_1 . By the splitting criterion of BUILDTOPDOWNDT and Fact 2.1, we have that

$$\text{Inf}_{x_1}(f_{\bar{\pi}}) \geq \text{Inf}_{y_j}(f_{\bar{\pi}}) \quad \text{for all } j \in [k].$$

Since

$$\text{Inf}_{x_1}(f_{\bar{\pi}}) \leq \text{Inf}_{x_1}(\text{TRIBES}_{\ell}) + \text{Inf}_{x_1}(\text{TRIBES}'_{\ell}) \leq O\left(\frac{\log \ell}{\ell}\right)$$

and

$$\text{Inf}_{y_j}(f_{\bar{\pi}}) = \Pr[\text{REST}_{\ell}(\mathbf{x}) = 1] \cdot \text{Inf}_{y_j}(\text{MAJ}_k(y)_{\bar{\pi}}) = 0.001 \cdot \text{Inf}_{y_j}(\text{MAJ}_k(y)_{\bar{\pi}}),$$

it follows that $\text{Inf}_{y_j}(\text{MAJ}_k(y)_{\bar{\pi}}) \leq O(\frac{\log \ell}{\ell})$. The claim follows by choosing c_1 to be a sufficiently small constant in Equation (2). \square

Corollary A.4. *There is a universal constant c_3 such that the following holds. Let $(\mathbf{x}, \mathbf{y}) \sim \{0, 1\}^k \times \{0, 1\}^{\ell}$ be a uniform random input, and $\pi_{(\mathbf{x}, \mathbf{y})}$ be the corresponding root-to-leaf path in T that (\mathbf{x}, \mathbf{y}) follows. The probability that $\pi_{(\mathbf{x}, \mathbf{y})}$ queries an x -variable before at least $c_3 k / \log k$ many y -variables is at most 0.001.*

Proof. Call an input (x, y) *bad* if $\pi_{(x, y)}$ queries an x -variable before $c_3 k$ many y -variables. Let $\bar{\pi}$ denote the truncation of $\pi_{(x, y)}$ to its prefix before the first query to an x -variable. By Claim A.3, we have that $\text{Inf}_{y_j}(\text{MAJ}_k(y)_{\bar{\pi}}) \leq \frac{1}{100\sqrt{k}}$, and so the discrepancy between the number of 0's and 1's in $\bar{\pi}$ must be $\Omega(\sqrt{k})$. Therefore, we can bound

$$\begin{aligned} \Pr[(\mathbf{x}, \mathbf{y}) \text{ is bad}] &\leq \sum_{t=1}^{c_3 k / \log k} \left(\Pr_{\mathbf{b} \sim \text{Bin}(t, \frac{1}{2})} \left[\left| \mathbf{b} - \frac{t}{2} \right| \geq \Omega(\sqrt{k}) \right] \right) \\ &\leq \sum_{t=1}^{c_3 k / \log k} \left(e^{-\Theta(k/t)} \right) \quad \text{(Hoeffding's inequality)} \\ &\leq c_3 k \cdot e^{-\Theta(\log k / c_3)} \leq o_k(1). \end{aligned}$$

where the final inequality holds by choosing c_3 to be a sufficiently small constant. \square

We are now ready to prove Lemma A.2. Let $\xi : \{0, 1\}^k \times \{0, 1\}^{\ell} \rightarrow \{0, 1\}$ be the indicator

$$\xi(x, y) := \mathbb{1} \left[\pi_{(x, y)} \text{ queries an } x\text{-variable before at least } \frac{c_3 k}{\log k} \text{ many } y\text{-variables} \right].$$

By Corollary A.4 we have that $\Pr[\xi(\mathbf{x}, \mathbf{y}) = 1] \leq 0.001$. If (x, y) is such that $\Pr[\xi(x, y) = 0]$, then either

1. $|\pi_{(x, y)}| \geq \frac{c_3 k}{\log k}$, or
2. $|\pi_{(x, y)}| < \frac{c_3 k}{\log k}$ and $\pi_{(x, y)}$ does not query any x -variables.

Since the fraction of inputs that follow any specific path of length $\geq c_3 k / \log k$ is at most $2^{-c_3 k / \log k}$, and the size of T is $2^{c_2 k / \log k}$ by assumption, choosing $c_2 = \frac{1}{2} c_3$ ensures that the fraction of inputs (x, y) such that $\pi_{(x, y)}$ falls into the first case above (i.e. $|\pi_{(x, y)}| \geq c_3 k / \log k$) is at most 0.001.

Therefore, at least a 0.998 fraction of inputs (x, y) are such that $\pi_{(x,y)}$ falls into the second case above. For such (x, y) 's, we have that

$$\begin{aligned}\Pr[f_{\pi_{(x,y)}} = 1] &\geq \Pr[\text{TRIBES}'_{\ell} = 1] = 0.499 \\ \Pr[f_{\pi_{(x,y)}} = 0] &\geq \Pr[\text{TRIBES}_{\ell} = 0] = 0.5,\end{aligned}$$

and so $\Pr[T_{\pi_{(x,y)}} \neq f_{\pi_{(x,y)}}] \geq 0.49$. We conclude that $\Pr[T(\mathbf{x}, \mathbf{y}) \neq f(\mathbf{x}, \mathbf{y})] \geq 0.998 \times 0.499 > 0.49$, which completes the proof of [Lemma A.2](#).

B. Proof of Theorem 1

We first prove [Theorem 5](#).

Proof of Theorem 5. First, without loss of generality, we can assume that \mathcal{D} is the uniform distribution on $[0, 1]^n$. Otherwise, we can transform each variable by its CDF. This also means that we will set $\theta^* = \frac{1}{2}$.

We will use [Theorem 4](#) as a black box in our proof of [Theorem 5](#). In order to do this, we will relate T , a decision tree with real-valued inputs, to a specially constructed decision tree with boolean-valued inputs. Each input of T will be encoded into a boolean vector of width $w = O(\log(s/\varepsilon))$ using the encoder $E : [0, 1] \rightarrow \{\pm 1\}^w$ defined as follows.

$$E(x) = \text{BINARY}(\lfloor x \cdot 2^w \rfloor),$$

where `BINARY` is the function that encodes an integer in binary form. We want to ensure that if two inputs are encoded as the same Boolean vector evaluate to the same value. This requires rounding the thresholds in T .

Lemma B.1. *Let $T : [0, 1]^n \rightarrow \{0, 1\}$ be a balanced size- s tree and $\varepsilon > 0$. Then for $w = O(\log(s/\varepsilon))$ let T_{round} be the tree formed by rounding all thresholds of T to the nearest 2^{-w} . Then $\text{dist}(T, T_{\text{round}}) \leq \varepsilon/2$.*

Proof of Lemma B.1. Since T and T_{round} have the same label at every leaf, if $T(x) \neq T_{\text{round}}(x)$, it means that x reaches different leaves in T and T_{round} . Fix some leaf ℓ of T . Then, ℓ is reached by some subcube of inputs of the form

$$a_i \leq x_i < b_i \quad \text{for } i = 1 \in [n]$$

The corresponding leaf, ℓ_{round} , of T_{round} is reached by all inputs satisfying the following, where $\text{round}(\cdot)$ rounds an input to the nearest multiple of 2^{-w} .

$$\text{round}(a_i) \leq x_i < \text{round}(b_i) \quad \text{for } i = 1 \in [n]$$

We will upper bound the probability that a randomly chosen $x \in [0, 1]^n$ reaches ℓ in T but does not reach ℓ_{round} in T_{round} . For that to occur, there must be one $i \in [n]$ such that $x_i \in [a_i, \text{round}(a_i)]$ or $x_i \in [\text{round}(b_i), b_i]$. Since $|\text{round}(z) - z| \leq 2^{-w}$, for any fixed i , x_i falls in one of those ranges with probability at most $2 \cdot 2^{-w}$.

Furthermore, since T is balanced, ℓ has depth at most $O(\log(s))$. This means that for all but $O(\log(s))$ choices for i , $a_i = 0$ and $b_i = 1$, in which case there is no chance that $x_i \in [a_i, \text{round}(a_i)]$ or $x_i \in [\text{round}(b_i), b_i]$. By union bounding the up to $O(\log(s))$ input coordinates that matter for ℓ , we have that the probability that x reaches ℓ in T but does not reach ℓ_{round} in T_{round} is at most $O(\log(s) \cdot 2^{-w})$.

By union bound over the s leaves, the probability x reaches a different leaf in T as in T_{round} is at most $O(s \log s \cdot 2^{-w})$. Setting this equal to $\varepsilon/2$ and solving for w yields the desired result. \square

We are now ready to complete the proof of [Theorem 5](#). Note that

$$\begin{aligned}\text{bias}(f) - \text{dist}(f, T_{\text{round}}) &\geq \text{bias}(f) - (\text{dist}(f, T) + \text{dist}(T, T_{\text{round}})) \\ &\geq \text{bias}(f) - \text{dist}(f, T) - \frac{\varepsilon}{2} \\ &\geq \frac{\varepsilon}{2}.\end{aligned}$$

We will create a tree with Boolean inputs, $S : \{\pm 1\}^{wn} \rightarrow \{0, 1\}$, satisfying the following relation

$$S(E(x_1), \dots, E(x_n)) = T_{\text{round}}(x) \quad \text{for all } x \in [0, 1]^n.$$

The decision tree computing S is the tree computing T where each threshold is replaced with the tree of size at most w specifying that threshold. Replacing each node of T with a tree of size w creates a tree that is a factor of w^d larger than T , where d is the depth of T . Using $w = O(\log(s/\varepsilon))$, we then have that S has size

$$O(\log(s/\varepsilon))^d = 2^{O(\log(s) \log \log(s/\varepsilon))}.$$

We are now ready to apply [Theorem 4](#). We have determined that $\text{bias}(f) - \text{dist}(f, T_{\text{round}}) \geq \varepsilon/2$, and that T_{round} is encodable as a decision tree of size $2^{\log(s) \log(\log(s/\varepsilon))}$ with Boolean inputs. By [Theorem 4](#), there is a bit of the encoding with influence, with respect to f , at least $\Omega(\varepsilon / \log(s) \log \log(s/\varepsilon))$.

All that remains is to show that the most influential bit of the encoding, is a single threshold of the form $\mathbb{1}[x_i \geq \frac{1}{2}]$. All bits of the encoding for a variable, x_i , restrict x_i to exactly half of the space in $[0, 1]$. Since f is monotone, the most influential such restriction is $\mathbb{1}[x_i \geq \frac{1}{2}]$. Hence, the split $\mathbb{1}[x_i \geq \frac{1}{2}]$ has influence at least $\Omega(\varepsilon / \log(s) \log \log(s/\varepsilon))$. Since f is monotone, this is also the correlation of $\mathbb{1}[x_i \geq \frac{1}{2}]$ with f , proving [Theorem 5](#). \square

The remainder of the proof of [Theorem 1](#) is the same as the proof of [Theorem 2](#). Let T° be the size- $(j+1)$ partial tree built BUILDTOPDOWNDT after j iterations. As long as $\text{dist}(f, T_f^\circ) > \text{balanced_opt}_s + \varepsilon$, we know there is a split that results in purity gain of at least⁶

$$\Omega\left(\frac{\kappa \cdot \varepsilon^2}{j \cdot (\log(s) \log \log(s/\varepsilon))^2}\right) \left(\right)$$

Therefore, after we have run for

$$\begin{aligned} t &= 2^{O(\log(s) \log \log(s/\varepsilon))^2 / \kappa \varepsilon^2} \\ &= s^{\tilde{O}((\log s) / \kappa \varepsilon^2)} \end{aligned}$$

we must have that $\text{dist}(f, T_f^\circ) < \text{balanced_opt}_s + \varepsilon$, proving [Theorem 1](#).

C. The Realizable Case

([Blanc et al., 2020](#))’s work on the realizable setting analyzed the performance of a *variant* of the top-down heuristics; their variant does not correspond to any impurity function \mathcal{G} . Consider BUILDTOPDOWNDT_{Inf}, defined in [Figure 2](#).

⁶From [Theorem 5](#), we actually know there is a split with this much purity that is the median of one input coordinate, although BUILDTOPDOWNDT may choose one that is not the median if it results in more purity gain.

BUILDTOPDOWNDT_{Inf}(f, t):

Initialize T° to be the empty tree.

while ($\text{size}(T^\circ) < t$ {

1. (Score) For every leaf ℓ in T° , let $x_{i(\ell)}$ denote the most influential variable of the subfunction f_ℓ :

$$\text{Inf}_{i(\ell)}(f_\ell) \geq \text{Inf}_j(f_\ell) \quad \text{for all } j \in [n].$$

Assign ℓ the score:

$$\text{score}_{\text{Inf}}(\ell) := 2^{-|\ell|} \cdot \text{Inf}_{i(\ell)}(f_\ell),$$

where $|\ell|$ denotes the depth of ℓ in T° .

2. (Split) Let ℓ^* be the leaf with the highest score. Grow T° by replacing ℓ^* with a query to $x_{i(\ell^*)}$.

}

Output f -completion of T° .

Figure 2. Top-down heuristic for building a decision tree approximation for f from (Blanc et al., 2020)

Lemma C.1 ((Theorem 5 of (Blanc et al., 2020))). *For every $\varepsilon \in (0, \frac{1}{2})$ and monotone function $f : \{\pm 1\}^n \rightarrow \{0, 1\}$ exactly computable by a size s decision tree,*

$$\text{TOPDOWNERROR}_{\text{Inf}}(f, s^{O(\sqrt{\log s}/\varepsilon)}) \leq \varepsilon.$$

We will show that as a consequence of Lemma C.1, a similar upper bound can be shown for BUILDTOPDOWNDT _{\mathcal{G}} for any strongly-concave \mathcal{G} . We first briefly summarize the proof of Lemma C.1 given in (Blanc et al., 2020). They define the following potential function:

$$u_f(T^\circ) := \sum_{\text{leaves } \ell \in T^\circ} 2^{-|\ell|} \cdot \text{Inf}(f_\ell).$$

For monotone functions $f : \{\pm 1\}^n \rightarrow \{0, 1\}$ that are computable by size- s decision trees, a result of (O’Donnell & Servedio, 2007) gives the following bound:

$$u_f(\text{empty tree}) = \text{Inf}(f) \leq \sqrt{\log s}.$$

Furthermore, they show that if T° is the size- $(j+1)$ tree built after j iterations of BUILDTOPDOWNDT_{Inf}(f), and $\text{dist}(f, T_f^\circ) \geq \varepsilon$, then the leaf, ℓ^* , selected in the next iteration satisfies.

$$\text{score}(\ell^*) \geq \frac{\varepsilon}{(j+1) \log s}.$$

Since $u_f(T^\circ)$ decreases by the score of the leaf selected, and $u_f(T^\circ)$ upper bounds $\text{dist}(f, T_f^\circ)$, (Blanc et al., 2020) are able to conclude that a growing a decision tree of size

$$2^{O(\sqrt{\log s} \cdot \log(s)/\varepsilon)} = s^{O(\sqrt{\log s}/\varepsilon)}$$

suffices to ensure an error of $\leq \varepsilon$. Using this same proof outline, we will establish the same guarantee for BUILDTOPDOWNDT _{\mathcal{G}} for any strongly concave impurity function \mathcal{G} :

Theorem 6 (Extending Theorem 5 of (Blanc et al., 2020) to actual top-down heuristics). *Let \mathcal{G} be any κ -strongly concave impurity function, $\varepsilon \in (0, \frac{1}{2})$, and $f : \{\pm 1\}^n \rightarrow \{0, 1\}$ be a monotone function computable by a size- s decision tree. Then,*

$$\text{TOPDOWNERROR}_{\mathcal{G}}(f, s^{O(\sqrt{\log s}/\varepsilon)})/\kappa \leq \varepsilon.$$

Proof. Recalling [Proposition 2.2](#), for any leaf $\ell \in T^\circ$, the variable with maximum influence in f_ℓ will also be the variable that results in the maximum purity gain when split. Therefore, $\text{BUILDTOPDOWNDT}_{\mathcal{G}}$ will always choose the same variable at any leaf to split as $\text{BUILDTOPDOWNDT}_{\text{Inf}}$, but may just choose a different order of leaves to split. Splitting extra leaves can only decrease the error, so once $\text{BUILDTOPDOWNDT}_{\mathcal{G}}$ has split every leaf that $\text{BUILDTOPDOWNDT}_{\text{Inf}}$ would in $s^{O(\sqrt{\log s}/\varepsilon)}$ iterations, the resulting tree must have error less than ε .

We know that running $\text{BUILDTOPDOWNDT}_{\text{Inf}}$ for $s^{O(\sqrt{\log s}/\varepsilon)}$ iterations is sufficient to ensure an error of at most ε . Let T° be the size- $(j+1)$ tree built after j iterations of $\text{BUILDTOPDOWNDT}_{\text{Inf}}(f)$. ([Blanc et al., 2020](#)) proved that if $\text{dist}(f, T_f^\circ) \geq \varepsilon$, then the leaf, ℓ^* , selected in the next iteration satisfies:

$$\text{score}(\ell^*) \geq \frac{\varepsilon}{(j+1) \log s}.$$

Therefore, we can substitute in $j = s^{O(\sqrt{\log s}/\varepsilon)}$ to find that if $\text{BUILDTOPDOWNDT}_{\text{Inf}}$ has split all leaves with score at least

$$\frac{\varepsilon}{s^{O(\sqrt{\log s}/\varepsilon)} \cdot \log s} = \frac{1}{s^{O(\sqrt{\log s}/\varepsilon)}},$$

the tree it has built must have error ε . We will show that $\text{BUILDTOPDOWNDT}_{\mathcal{G}}$ will not take too long to split any leaf with score at least $1/s^{O(\sqrt{\log s}/\varepsilon)}$, therefore proving an upper bound on the number of iterations it needs to reach error ε . Let ℓ be any leaf with score at least $1/s^{O(\sqrt{\log s}/\varepsilon)}$ and i be the index of its most influential variable. Then,

$$\begin{aligned} \mathcal{G}\text{-purity-gain}_f(T^\circ, \ell, x_i) &\geq 2^{|\ell|} \cdot \frac{\kappa}{32} \cdot \text{Inf}_i(f)^2 && \text{(Lemma 2.4)} \\ &\geq \frac{\kappa}{32} \cdot (2^{|\ell|} \cdot \text{Inf}_i(f))^2 \\ &= \frac{\kappa}{32} \cdot \text{score}(\ell)^2 \\ &\geq \frac{\kappa}{s^{O(\sqrt{\log s}/\varepsilon)}}. \end{aligned}$$

Let ℓ^* be some leaf with larger purity gain than ℓ , and i^* be the associated variable that is split. Then, since $\mathcal{G}\text{-purity-gain}_f(T^\circ, \ell^*, x_{i^*}) \leq 2^{-|\ell^*|}$, we must have that $|\ell^*| < \log(s^{O(\sqrt{\log s}/\varepsilon)}/\kappa)$. There are at most $s^{O(\sqrt{\log s}/\varepsilon)}/\kappa$ possible such nodes, so after constructing a tree of size $s^{O(\sqrt{\log s}/\varepsilon)}/\kappa$, $\text{BUILDTOPDOWNDT}_{\mathcal{G}}$ must split ℓ , and we conclude that it must have achieved error at most ε . \square