

A Other related work

Kearns and Mansour [KM99] (see also [Kea96]) analyzed top-down impurity-based heuristics from the perspective of *boosting*, where the attributes queried in the tree are viewed as weak hypotheses.

Recent work of Blanc et al. [BLT20b] gives a top-down algorithm for learning decision trees that achieves provable guarantees for all target functions f . However, their algorithm makes crucial use of *membership queries*, which significantly limits its practical applicability and relevance. Furthermore, their guarantees only hold in the realizable setting, requiring that f is itself a size- s decision tree (i.e. $\text{opt}_s = 0$).

There has been extensive work in the learning theory literature on learning the concept class of decision trees [EH89, Blu92, KM93, OS07, GKK08, HKY18, CM19]. However, none of these algorithms proceed in a top-down manner like the practical heuristics that are the focus of this work; indeed, with the exception [EH89], these algorithms do not return a decision tree as their hypothesis. ([EH89]'s algorithm constructs its decision tree hypothesis in a *bottom-up* manner.)

B Proof of Fact 2.1

Fact 2.1 is a simple consequence of the following lemma, whose proof also appears in [Jon16]:

Lemma B.1. *For all $f : \{\pm 1\}^n \rightarrow \{\pm 1\}$ and $i \in [n]$,*

$$\text{NS}_\delta(f) = \frac{1}{2} \text{NS}_\delta(f_{x_i=-1}) + \frac{1}{2} \text{NS}_\delta(f_{x_i=1}) + \frac{\delta}{2(1-\delta)} \cdot \text{Inf}_i^{(\delta)}(f).$$

Proof. Let $\mathbf{x} \sim \{\pm 1\}^n$ be uniform random, and $\tilde{\mathbf{x}} \sim_\delta \mathbf{x}$ be a δ -noisy copy of \mathbf{x} . We first note that

$$\begin{aligned} \mathbb{E}[f(\mathbf{x})f(\tilde{\mathbf{x}})] &= \Pr[\mathbf{x}_i = \tilde{\mathbf{x}}_i] \cdot \mathbb{E}[f(\mathbf{x})f(\tilde{\mathbf{x}}) \mid \mathbf{x}_i = \tilde{\mathbf{x}}_i] + \Pr[\mathbf{x}_i \neq \tilde{\mathbf{x}}_i] \cdot \mathbb{E}[f(\mathbf{x})f(\tilde{\mathbf{x}}) \mid \mathbf{x}_i \neq \tilde{\mathbf{x}}_i] \\ &= \left(1 - \frac{\delta}{2}\right) \left(\frac{1}{2} \mathbb{E}[f(\mathbf{x}^{i=1})f(\tilde{\mathbf{x}}^{i=1})] + \frac{1}{2} \mathbb{E}[f(\mathbf{x}^{i=-1})f(\tilde{\mathbf{x}}^{i=-1})]\right) \\ &\quad + \frac{\delta}{2} \left(\frac{1}{2} \mathbb{E}[f(\mathbf{x}^{i=1})f(\tilde{\mathbf{x}}^{i=-1})] + \frac{1}{2} \mathbb{E}[f(\mathbf{x}^{i=-1})f(\tilde{\mathbf{x}}^{i=1})]\right). \end{aligned} \quad (7)$$

Next, we have that

$$\begin{aligned} \mathbb{E}[D_i f(\mathbf{x})D_i f(\tilde{\mathbf{x}})] &= \frac{1}{4} \mathbb{E}[(f(\mathbf{x}^{i=1}) - f(\mathbf{x}^{i=-1}))(f(\tilde{\mathbf{x}}^{i=1}) - f(\tilde{\mathbf{x}}^{i=-1}))] \\ &= \frac{1}{4} \mathbb{E}[f(\mathbf{x}^{i=1})f(\tilde{\mathbf{x}}^{i=1})] + \frac{1}{4} \mathbb{E}[f(\mathbf{x}^{i=-1})f(\tilde{\mathbf{x}}^{i=-1})] \\ &\quad - \frac{1}{4} \mathbb{E}[f(\mathbf{x}^{i=1})f(\tilde{\mathbf{x}}^{i=-1})] - \frac{1}{4} \mathbb{E}[f(\mathbf{x}^{i=-1})f(\tilde{\mathbf{x}}^{i=1})]. \end{aligned} \quad (8)$$

Combining Equations (7) and (8),

$$\begin{aligned} \mathbb{E}[f(\mathbf{x})f(\tilde{\mathbf{x}})] &= \frac{1}{2} \mathbb{E}[f(\mathbf{x}^{i=1})f(\tilde{\mathbf{x}}^{i=1})] + \frac{1}{2} \mathbb{E}[f(\mathbf{x}^{i=-1})f(\tilde{\mathbf{x}}^{i=-1})] - \delta \mathbb{E}[D_i f(\mathbf{x})D_i f(\tilde{\mathbf{x}})] \\ &= \frac{1}{2} \mathbb{E}[f_{x_i=1}(\mathbf{x})f_{x_i=1}(\tilde{\mathbf{x}})] + \frac{1}{2} \mathbb{E}[f_{x_i=-1}(\mathbf{x})f_{x_i=-1}(\tilde{\mathbf{x}})] - \frac{\delta}{1-\delta} \cdot \text{Inf}_i^{(\delta)}(f). \end{aligned}$$

Since $\text{NS}_\delta(f) = \Pr[f(\mathbf{x}) \neq f(\tilde{\mathbf{x}})] = \frac{1}{2} - \frac{1}{2} \mathbb{E}[f(\mathbf{x})f(\tilde{\mathbf{x}})]$, the lemma follows from the above by rearranging. \square

Proof of Fact 2.1. We first note that

$$\begin{aligned} \text{NS}_\delta(f, T_{\ell^* \rightarrow x_i}^\circ) &= \sum_{\text{leaves } \ell \in T_{\ell^* \rightarrow x_i}^\circ} 2^{-|\ell|} \cdot \text{NS}_\delta(f_\ell) \\ &= \sum_{\text{leaves } \ell \in T^\circ} 2^{-|\ell|} \cdot \text{NS}_\delta(f_\ell) \\ &\quad + 2^{-(|\ell^*|+1)} \cdot \text{NS}_\delta(f_{\ell^*, x_i=-1}) + 2^{-(|\ell^*|+1)} \cdot \text{NS}_\delta(f_{\ell^*, x_i=1}) - 2^{-|\ell^*|} \cdot \text{NS}_\delta(f_{\ell^*}) \\ &= \text{NS}_\delta(f, T^\circ) + 2^{-|\ell^*|} \left(\frac{1}{2} \text{NS}_\delta(f_{\ell^*, x_i=-1}) + \frac{1}{2} \text{NS}_\delta(f_{\ell^*, x_i=1}) - \text{NS}_\delta(f_{\ell^*})\right). \end{aligned}$$

Applying Lemma B.1 with its ' f ' being f_{ℓ^*} , we have that

$$\frac{1}{2} \text{NS}_\delta(f_{\ell^*, x_i=-1}) + \frac{1}{2} \text{NS}_\delta(f_{\ell^*, x_i=1}) - \text{NS}_\delta(f_{\ell^*}) = -\frac{\delta}{2(1-\delta)} \cdot \text{Inf}_i^{(\delta)}(f_{\ell^*}),$$

and this completes the proof. \square

C Proof of Theorem 3

Proof. We apply Theorem 2 with ‘ T ’ being T^* , ‘ g ’ being $f_{\delta}^{\leq d}$, and ρ being the semimetric $\rho(a, b) = (a - b)^2/2$. As shown by [OSSS05] (and as can be easily verified), $\text{Def}_k(\rho) \leq k$ for this choice of ρ , and so

$$\text{Cov}_{\rho}(T^*, f_{\delta}^{\leq d}) \leq k \sum_{i=1}^n \lambda_i(T^*) \cdot \frac{1}{2} \mathbb{E} [(f_{\delta}^{\leq d}(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}^{\sim i}))^2]. \quad (9)$$

We first analyze the quantity on the LHS of Equation (9). For $\mathbf{x}, \mathbf{x}' \sim \{\pm 1\}^n$ uniform and independent,

$$\begin{aligned} \text{Cov}_{\rho}(T^*, f_{\delta}^{\leq d}) &= \frac{1}{2} (\mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}'))^2] - \mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}))^2]) \\ &\geq \frac{1}{4} \mathbb{E} [(f_{\delta}^{\leq d}(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}'))^2] - \mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}))^2] \\ &= \frac{1}{2} \text{Var}(f_{\delta}^{\leq d}) - \mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}))^2], \end{aligned} \quad (10)$$

where the inequality uses the ‘‘almost-triangle’’ inequality $(a - c)^2 \leq 2((a - b)^2 + (b - c)^2)$ for $a, b, c \in \mathbb{R}$. Furthermore, we have

$$\begin{aligned} \text{Var}(f_{\delta}) &= \sum_{S \neq \emptyset} (1 - \delta)^{2|S|} \widehat{f}(S)^2 && \text{(Fourier formulas for } f_{\delta} \text{ (5) and variance (4))} \\ &= \sum_{\substack{S \neq \emptyset \\ |S| \leq d}} (1 - \delta)^{2|S|} \widehat{f}(S)^2 + \sum_{|S| > d} (1 - \delta)^{2|S|} \widehat{f}(S)^2 \\ &\leq \text{Var}(f_{\delta}^{\leq d}) + \sum_{|S| > d} e^{(-\delta)2|S|} \widehat{f}(S)^2 && (1 + a \leq e^a) \\ &\leq \text{Var}(f_{\delta}^{\leq d}) + e^{-2d\delta} \sum_{|S| > d} \widehat{f}(S)^2 && \text{(Since } |S| > d) \\ &\leq \text{Var}(f_{\delta}^{\leq d}) + e^{-2d\delta} && \text{(Parseval's identity (3): } \sum_{S \subseteq [n]} \widehat{f}(S)^2 = 1) \\ &\leq \text{Var}(f_{\delta}^{\leq d}) + \varepsilon. && \text{(Since } d = \log(1/\varepsilon)/\delta) \end{aligned}$$

Similarly,

$$\begin{aligned} \mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}))^2] &\leq 2 (\mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}(\mathbf{x}))^2] + \mathbb{E} [(f_{\delta}(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}))^2]) \\ &&& \text{('‘almost-triangle’’ inequality)} \\ &= 2 \left(\mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}(\mathbf{x}))^2] + \sum_{|S| > d} (1 - \delta)^{|S|} \widehat{f}(S)^2 \right) \\ &\leq 2 (\mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}(\mathbf{x}))^2] + \varepsilon). && \text{(Since } d = \log(1/\varepsilon)/\delta) \end{aligned}$$

Combining these bounds with Equation (10), we have the following lower bound on the LHS of Equation (9):

$$\begin{aligned} \text{Cov}(T^*, f_{\delta}^{\leq d}) &\geq \frac{1}{2} (\text{Var}(f_{\delta}) - \varepsilon) - (2 \mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}(\mathbf{x}))^2] + 2\varepsilon). \\ &= \frac{1}{2} \text{Var}(f_{\delta}) - (2 \mathbb{E} [(T^*(\mathbf{x}) - f_{\delta}(\mathbf{x}))^2] + \frac{5}{2}\varepsilon). \end{aligned} \quad (11)$$

We now turn to analyzing the RHS of Equation (9):

$$\begin{aligned}
& k \sum_{i=1}^n \lambda_i(T^*) \cdot \frac{1}{2} \mathbb{E} [(f_{\delta}^{\leq d}(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}^{\sim i}))^2] \\
&= k \sum_{i=1}^n \lambda_i(T^*) \cdot \frac{1}{4} \mathbb{E} [(f_{\delta}^{\leq d}(\mathbf{x}) - f_{\delta}^{\leq d}(\mathbf{x}^{\oplus i}))^2] \quad (\mathbf{x}^{\oplus i} = \mathbf{x} \text{ with its } i\text{-th coordinate flipped}) \\
&= k \sum_{i=1}^n \lambda_i(T^*) \cdot \mathbb{E} [D_i f_{\delta}^{\leq d}(\mathbf{x})^2] \\
&= k \sum_{i=1}^n \lambda_i(T^*) \cdot \text{Inf}_i(f_{\delta}^{\leq d}) \quad (\text{Definition 2}) \\
&= k \cdot \max_{i \in [n]} \left\{ \text{Inf}_i(f_{\delta}^{\leq d}) \right\} \cdot \sum_{i=1}^n \lambda_i(T^*) \leq k \cdot \max_{i \in [n]} \left\{ \text{Inf}_i(f_{\delta}^{\leq d}) \right\} \cdot \log s, \quad (12)
\end{aligned}$$

where the final inequality holds because

$$\sum_{i=1}^n \lambda_i(T^*) = \sum_{i=1}^n \Pr[T^* \text{ queries } \mathbf{x}_i] = \sum_{\text{leaves } \ell \in T^*} 2^{-|\ell|} \cdot |\ell| \leq \log s.$$

Finally, we note that:

$$\begin{aligned}
\text{Inf}_i(f_{\delta}^{\leq d}) &= \sum_{\substack{S \ni i \\ |S| \leq d}} (1 - \delta)^{2|S|} \widehat{f}(S)^2 \quad (\text{Fourier formula for influence; Definition 2}) \\
&\leq \sum_{\substack{S \ni i \\ |S| \leq d}} (1 - \delta)^{|S|} \widehat{f}(S)^2 = \text{Inf}_i^{(\delta, d)}(f).
\end{aligned}$$

Combining this with Equations (9), (11) and (12) and rearranging completes the proof. \square

D Proofs of Facts 4.1 and 4.2 and Propositions E.1 and E.2

Proof of Fact 4.1. This follows by combining the bounds $\text{Inf}(T) \leq \log s$ (see e.g. [OS07]) and $\text{NS}_{\delta}(f) \leq \delta \cdot \text{Inf}(f)$ for all $f : \{\pm 1\}^n \rightarrow \{\pm 1\}$ [O'D14, Exercise 2.42]. \square

Proof of Fact 4.2. Let $\mathbf{x} \sim \{\pm 1\}^n$ be uniform random and $\tilde{\mathbf{x}} \sim_{\delta} \mathbf{x}$ be a δ -noisy copy of \mathbf{x} . Then

$$\begin{aligned}
\text{NS}_{\delta}(f) &= \Pr[f(\mathbf{x}) \neq f(\tilde{\mathbf{x}})] \\
&\leq \Pr[T(\mathbf{x}) \neq T(\tilde{\mathbf{x}})] + \Pr[T(\mathbf{x}) \neq f(\mathbf{x})] + \Pr[T(\tilde{\mathbf{x}}) \neq f(\tilde{\mathbf{x}})] \\
&\leq \text{NS}_{\delta}(T) + 2 \Pr[T(\mathbf{x}) \neq f(\mathbf{x})],
\end{aligned}$$

where the final inequality uses that fact that \mathbf{x} and $\tilde{\mathbf{x}}$ are distributed identically. \square

E The case analysis in the proof of Theorem 4

Case 1: $\mathbb{E}_{\ell}[\text{Var}((f_{\ell})_{\delta})] \geq 4 \mathbb{E}_{\ell}[\|(f_{\ell})_{\delta} - T_{\text{opt}}^{\text{trunc}}\|_2^2] + 7\varepsilon.$

In this case we claim that there is a leaf ℓ^* of T° with a high score, where we recall that the score of a leaf ℓ is defined to be

$$\text{score}(\ell) := 2^{-|\ell|} \cdot \max_{i \in [n]} \left\{ \text{Inf}_i^{(\delta, d)}(f_{\ell}) \right\}.$$

Applying Theorem 3 with its ‘ T^* ’ being $T_{\text{opt}}^{\text{trunc}}$ and its ‘ f ’ being f_ℓ for each leaf $\ell \in T^\circ$, we have that

$$\begin{aligned} \mathbb{E} \left[\max_{i \in [n]} \{ \text{Inf}_i^{(\delta, d)}(f_\ell) \} \right] &\geq \frac{\frac{1}{2} \mathbb{E}[\text{Var}(f_\ell)_\delta] - \left(2 \mathbb{E} \left[\|T_{\text{opt}}^{\text{trunc}} - (f_\ell)_\delta\|_2^2 \right] + \frac{5}{2} \varepsilon \right)}{\log(s/\varepsilon) \log s} \quad (\text{Theorem 3}) \\ &\geq \frac{\varepsilon}{\log(s/\varepsilon) \log s}, \end{aligned} \quad (13)$$

where the second inequality is by the assumption that we are in Case 1. Equivalently,

$$\sum_{\ell \in T^\circ} 2^{-|\ell|} \cdot \max_{i \in [n]} \{ \text{Inf}_i^{\delta, d}(f_\ell) \} \geq \frac{\varepsilon}{\log(s/\varepsilon) \log s},$$

and so there must exist a leaf $\ell^* \in T^\circ$ such that

$$\text{score}(\ell^*) = 2^{-|\ell^*|} \cdot \max_{i \in [n]} \{ \text{Inf}_i^{(\delta, d)}(f_{\ell^*}) \} \geq \frac{\varepsilon}{|T^\circ| \log(s/\varepsilon) \log s},$$

where $|T^\circ|$ denotes the size of T° .

Case 2: $\mathbb{E}[\text{Var}((f_\ell)_\delta)] < 4 \mathbb{E}[\|(f_\ell)_\delta - T_{\text{opt}}^{\text{trunc}}\|_2^2] + 7\varepsilon$.

In this case, we claim that $\text{error}_f(T_f^\circ) \leq O(\text{opt}_s + \kappa + \varepsilon)$. We will need a couple of propositions:

Proposition E.1. $\mathbb{E}[\|(f_\ell)_\delta - f_\ell\|_2^2] \leq 4\kappa$.

Proof. We first note that

$$\begin{aligned} \mathbb{E}_\ell [\|(f_\ell)_\delta - f_\ell\|_2^2] &\leq 2 \mathbb{E}_\ell [\|(f_\ell)_\delta - f_\ell\|_1] \quad (\text{Since } f_\ell \text{ and } (f_\ell)_\delta \text{ are } [-1, 1]\text{-valued}) \\ &= 2 \mathbb{E}_\ell \left[\mathbb{E}_{\mathbf{x}} [|(f_\ell)_\delta(\mathbf{x}) - f_\ell(\mathbf{x})|] \right] \\ &= 2 \mathbb{E}_\ell \left[\mathbb{E}_{\substack{\mathbf{x} \\ \tilde{\mathbf{x}} \sim_\delta \mathbf{x}}} [|(f_\ell)(\tilde{\mathbf{x}}) - f_\ell(\mathbf{x})|] \right] \\ &= 2 \mathbb{E}_\ell \left[2 \Pr_{\substack{\mathbf{x} \\ \tilde{\mathbf{x}} \sim_\delta \mathbf{x}}} [f_\ell(\tilde{\mathbf{x}}) \neq f_\ell(\mathbf{x})] \right] \\ &= 4 \mathbb{E}_\ell [\text{NS}_\delta(f_\ell)] = 4 \text{NS}_\delta(f, T^\circ). \end{aligned}$$

By Fact 2.1, we have that $\text{NS}_\delta(f, T^\circ) \leq \text{NS}_\delta(f)$, and the claim follows. \square

Proposition E.2. For each leaf $\ell \in T^\circ$, we have $\mathbb{E}[(f_\ell(\mathbf{x}) - \text{sign}(\mathbb{E}[f_\ell]))^2] \leq 2 \mathbb{E}[(f_\ell(\mathbf{x}) - c)^2]$ for all constants $c \in \mathbb{R}$.

Proof. Let $p := \Pr[f_\ell(\mathbf{x}) = 1]$ and assume without loss of generality that $p \geq \frac{1}{2}$. On one hand, we have that $\mathbb{E}[(f_\ell(\mathbf{x}) - \text{sign}(\mathbb{E}[f_\ell]))^2] = \mathbb{E}[(f_\ell(\mathbf{x}) - 1)^2] = 4(1 - p)$. On the other hand, since

$$\mathbb{E}[(f_\ell(\mathbf{x}) - c)^2] = p(1 - c)^2 + (1 - p)(1 + c)^2$$

this quantity is minimized for $c = 2p - 1$ and attains value $4p(1 - p)$ at this minimum. Therefore indeed

$$\min_c \{ \mathbb{E}[(f_\ell(\mathbf{x}) - c)^2] \} = 4p(1 - p) \geq 2(1 - p) = \frac{1}{2} \mathbb{E}[(f_\ell(\mathbf{x}) - \text{sign}(\mathbb{E}[f_\ell]))^2]$$

and the proposition follows. \square

With Propositions E.1 and E.2 in hand, we are ready to bound $\text{error}_f(T_f^\circ)$. Recall that T_f° is the completion of T° that we obtain by labeling each leaf ℓ with $\text{sign}(\mathbb{E}[f_\ell])$. Therefore,

$$\begin{aligned}
\text{error}_f(T_f^\circ) &= \mathbb{E}_\ell [\text{dist}(f_\ell, \text{sign}(\mathbb{E}[f_\ell]))] \\
&= \frac{1}{4} \mathbb{E}_\ell [\|f_\ell - \text{sign}(\mathbb{E}[f_\ell])\|_2^2] \\
&\leq \frac{1}{2} \mathbb{E}_\ell [\|f_\ell - \mathbb{E}[(f_\ell)_\delta]\|_2^2] && \text{(Proposition E.2)} \\
&\leq \mathbb{E}_\ell [\|f_\ell - (f_\ell)_\delta\|_2^2] + \mathbb{E}_\ell [\|(f_\ell)_\delta - \mathbb{E}[(f_\ell)_\delta]\|_2^2] \\
&\leq 4\kappa + \mathbb{E}_\ell [\text{Var}((f_\ell)_\delta)] && \text{(Proposition E.1)}
\end{aligned}$$

By the assumption that we are in Case 2, we have that:

$$\begin{aligned}
\mathbb{E}_\ell [\text{Var}((f_\ell)_\delta)] &< 4 \mathbb{E}_\ell [\|(f_\ell)_\delta - T_{\text{opt}}^{\text{trunc}}\|_2^2] + 7\varepsilon \\
&\leq 8 \left(4\kappa + \mathbb{E}_\ell [\|f_\ell - T_{\text{opt}}^{\text{trunc}}\|_2^2] \right) + 7\varepsilon && \text{(Proposition E.1)} \\
&= 8 \left(4\kappa + 4 \mathbb{E}_\ell [\text{dist}(f_\ell, T_{\text{opt}}^{\text{trunc}})] \right) + 7\varepsilon \\
&= 8(4\kappa + 4(\text{opt}_s + \varepsilon)) + 7\varepsilon \\
&\leq O(\text{opt}_s + \kappa + \varepsilon).
\end{aligned}$$

and so we have shown that $\text{error}_f(T_f^\circ) \leq O(\text{opt}_s + \kappa + \varepsilon)$.