

On the Consistency Rate of Decision Tree Learning Algorithms

Qin-Cheng Zheng, Shen-Huan Lyu, Shao-Qun Zhang, Yuan Jiang, Zhi-Hua Zhou
{ zhengqc, lvsh, zhangsq, jiangy, zhouzh }@lamda.nju.edu.cn

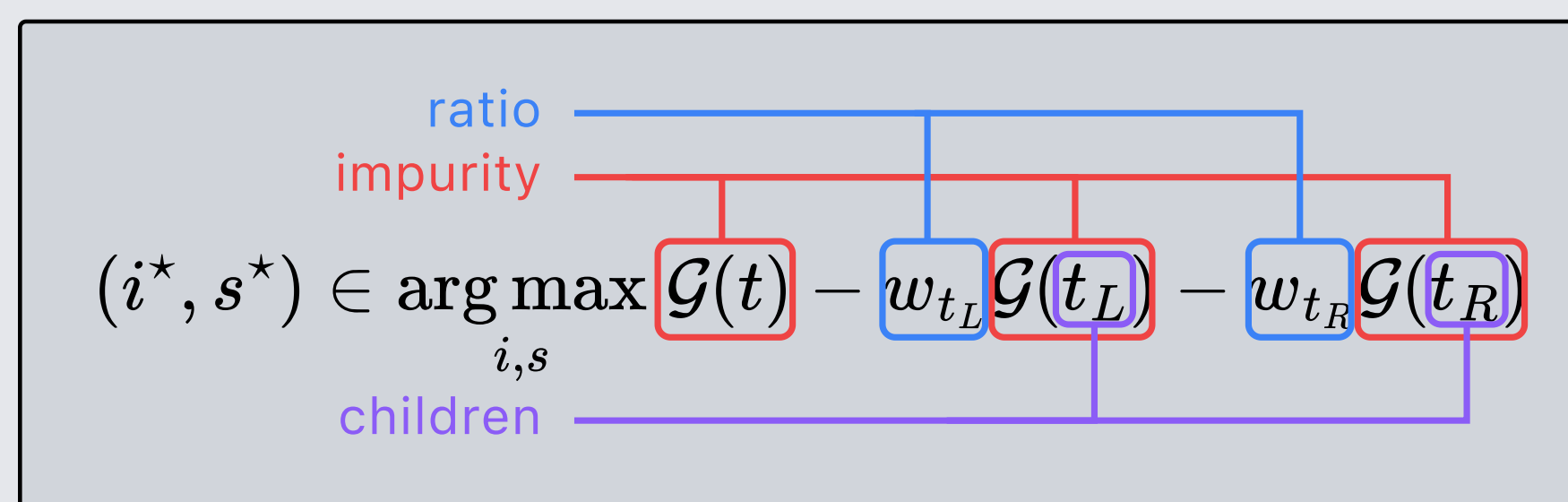


Paper in a Nutshell

1. We found that the **worst-case zero purity gain** leads to a serious obstacle for consistency analysis of CART.
2. We found that using Influence as the impurity measure can **always get a positive purity gain**, but an Influence oracle is required for generating a tree.
3. We propose the GridCART, an Influence-based CART, which not only can run **practically** but also is **consistent** with order $O(n^{-1/(d+2)})$.

Background

Heuristical algorithms maximize the **purity gain** to split each node to two children.

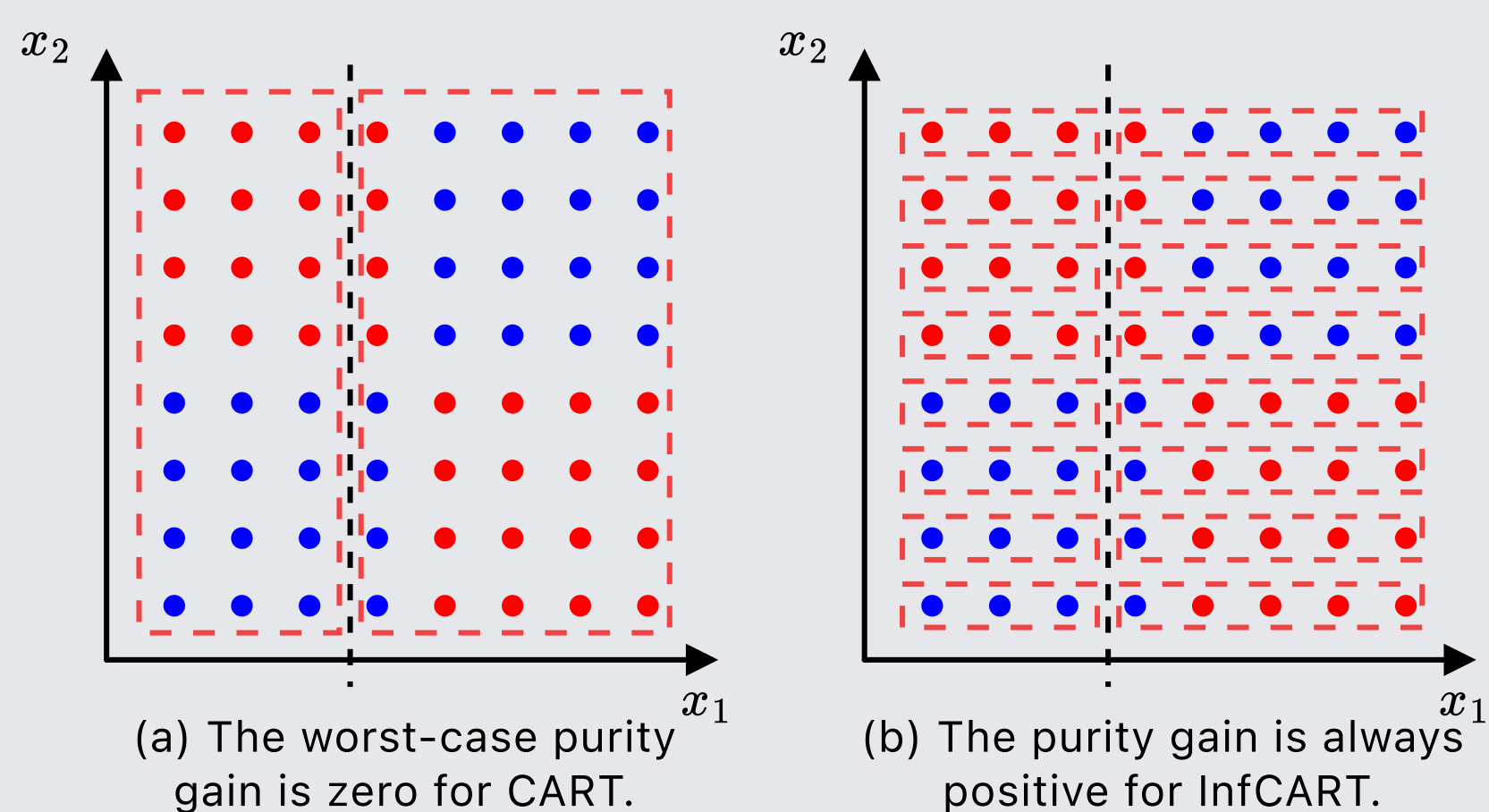


This type of algorithm succeeds in many real-world tasks, but the consistency is **far from clear**.

CART vs InfCART

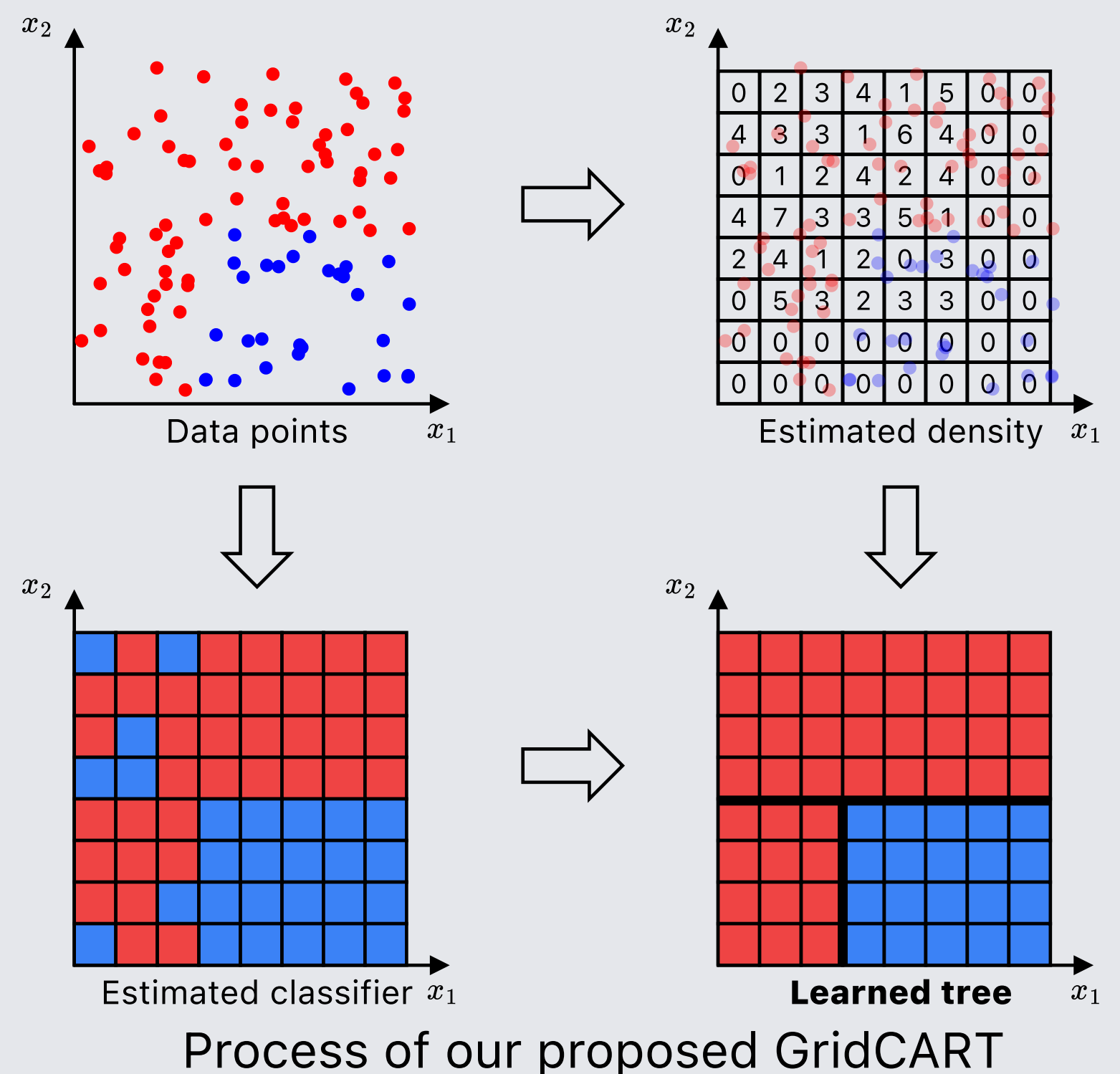
Algorithm	Impurity Measure
CART	$\mathcal{G}(t) \triangleq G\left(\mathbb{E}_{\mathbf{x}(j \neq i)} [\mathbb{E}_{\mathbf{x}(i)} [Y \mathbf{x} \in t]]\right)$
InfCART	$\mathcal{G}(t) \triangleq \mathbb{E}_{\mathbf{x}(j \neq i)} G(\mathbb{E}_{\mathbf{x}(i)} [Y \mathbf{x} \in t])$

CART always gets **zero purity gain** in the worst case see (a), while InfCART (a variant) calculates the average purity gain on every line and gets **a positive average purity gain** see (b).



Our Proposed GridCART

GridCART can **not only be practical but also always obtains a positive purity gain**. The key is to estimate the probability density, and further the Influence.



Theoretical Guarantees

Consistency rate for GridCART

Theorem 8 (informal). Under certain assumptions, the expected excess error of tree T_{K_n} learned by GridCART has the following upper bound

$$\mathbb{E}_{D_n} [R(T_{K_n})] - R^* \leq \mathcal{O} \left(\frac{1}{\text{depth}} \left(\frac{1}{K_n h_n^3} + \sqrt{\frac{1}{n h_n^d}} \right) \right)$$

cube size
sample size

Choosing $h_n = \Theta(n^{-1/(d+2)})$, $K_n = \Omega(n^{4/(d+2)})$, we obtain a consistency rate of order $\mathcal{O}(n^{-1/(d+2)})$.

GridCART can reach a consistency rate of order $\mathcal{O}(n^{-1/(d+2)})$ under certain assumptions. However, the consistency rate of CART is far from clear.

Algorithm	Advantage	Disadvantage
CART	Empirically practical	Worst-case zero purity gain
InfCART	Positive average purity gain	An Influence oracle required