

RHALE: Robust and Heterogeneity-aware Accumulated Local Effects

Vasilis Gkolemis^{1, 2} Theodore Dalamagas¹ Eirini Ntoutsis³ Christos Diou²

¹ATHENA Research Center

²Harokopio University of Athens

³Bundeswehr University of Munich

TL;DR

RHALE: Robust and Heterogeneity-aware ALE

- Robust: auto-bin splitting
- Heterogeneity: \pm from the average

keywords: eXplainable AI, ALE, Heterogeneity, Feature Effect

Motivation - ALE limitations

ALE (Apley2020) is a SotA feature effect method but it has two limitations:

- it does not quantify the heterogeneity, i.e., deviation of the instance-level effects from the main (average) effect
- it is vulnerable to poor approximations, due to the bin-splitting step

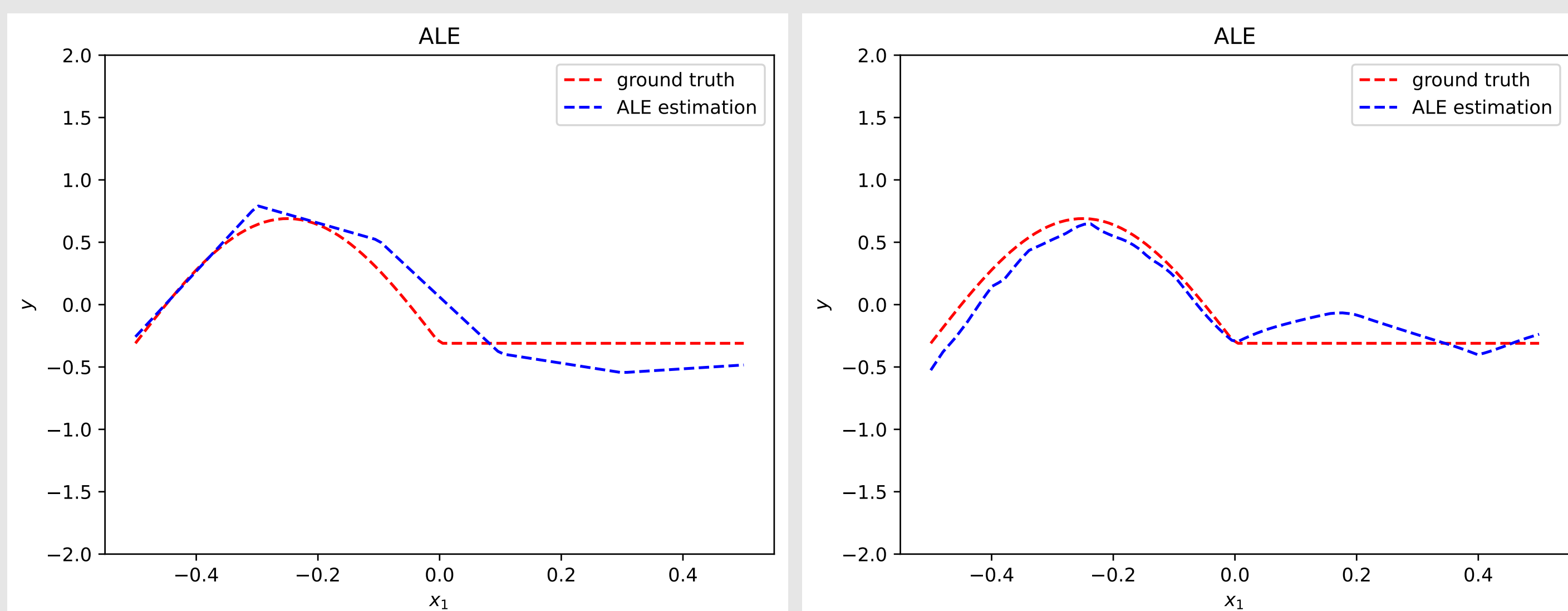


Figure 1. Left: ALE approximation with a narrow bin-splitting (5 bins). Right: ALE approximation with a dense bin-splitting (50 bins)

- Both approximations are bad:
 - Narrow bin-splitting hides fine-grain details
 - Dense bin-splitting is noisy (low samples per bin rate)
- No information regarding the heterogeneity

Simple approach: ALE + Heterogeneity

ALE main effect definition

$$f^{\text{ALE}}(x_s) = \int_{x_{s,\min}}^{x_s} \underbrace{\mathbb{E}_{X_c|X_s=z} [f^s(z, X_c)]}_{\mu(z)} \partial z$$

ALE main effect approximation

$$\hat{f}^{\text{ALE}}(x_s) = \Delta x \sum_k \frac{1}{|\mathcal{S}_k|} \sum_{i: \mathbf{x}^i \in \mathcal{S}_k} \underbrace{\left[\frac{\partial f}{\partial x_s}(x_s^i, \mathbf{x}_c^i) \right]}_{\text{bin effect: } \hat{\mu}(z)}$$

ALE heterogeneity definition

$$\sigma(x_s) = \sqrt{\int_{x_{s,\min}}^{x_s} \underbrace{\mathbb{E}_{X_c|X_s=z} [(f^s(z, X_c) - \mu(z))^2]}_{\sigma^2(z)} \partial z}$$

ALE heterogeneity approximation

$$\text{STD}(x_s) = \sqrt{\sum_{k=1}^{k_x} (z_k - z_{k-1})^2 \frac{1}{|\mathcal{S}_k| - 1} \sum_{i: \mathbf{x}^i \in \mathcal{S}_k} (f^s(\mathbf{x}^i) - \hat{\mu}(z_1, z_2))^2}$$

Simple but wrong: ALE + Heterogeneity

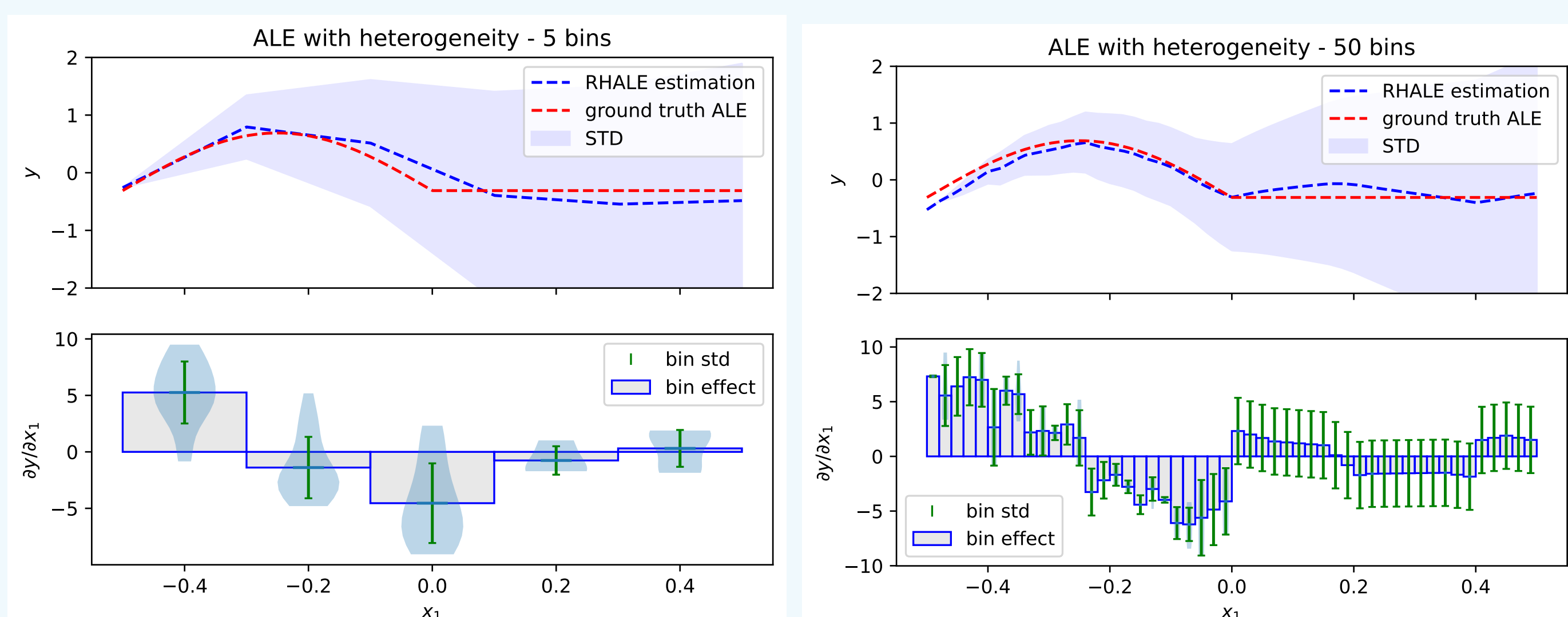


Figure 2. Left: approximation with narrow bin-splitting (5 bins) and (Right) with dense-bin splitting

- Fixed-size bin splitting can ruin the estimation of the heterogeneity

RHALE: Robust and Heterogeneity-aware ALE

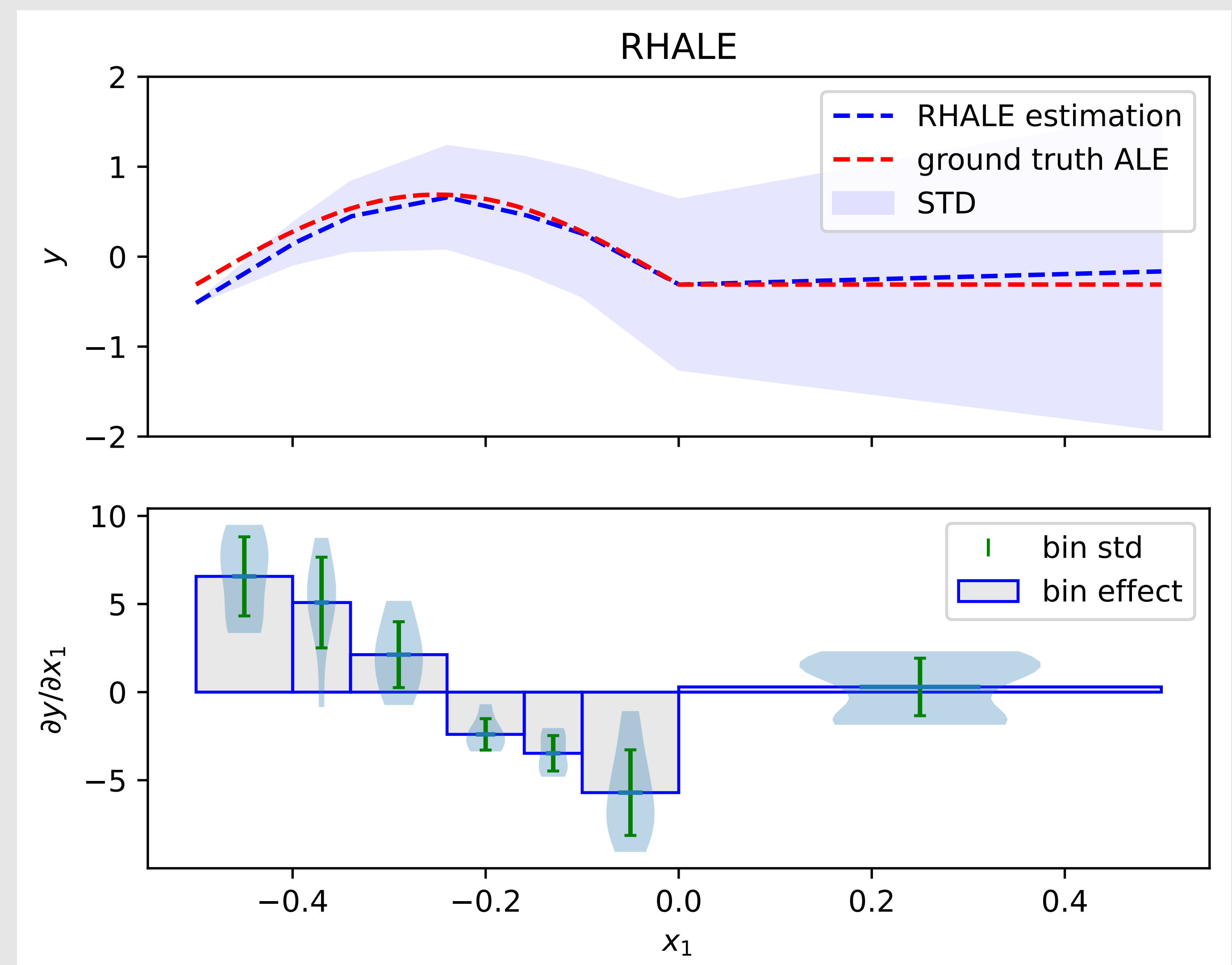


Figure 3. Narrow bins ($K = 40$) \Rightarrow limited $\frac{\text{samples}}{\text{bin}} \Rightarrow$ both plots are noisy

Simple but correct:

- Automatically finds the **optimal** bin-splitting
- Optimal \Rightarrow best approximation of the average (ALE) effect
- Optimal \Rightarrow best approximation of the heterogeneity

Optimal bin-splitting

In the paper, we **formally prove**:

- 1 the conditions under which a bin does not hide the finer details of the main effect
- 2 the conditions under which a bin is an unbiased approximator of the heterogeneity
- 3 that given (1) and (2), increasing bin size leads to a reduction in estimation variance

Based on the above, we formulate bin-splitting as an optimization problem and propose an efficient solution using dynamic programming.

Conclusion

In case you work with a differentiable model, as in Deep Learning, use RHALE to:

- quantify the heterogeneity of the ALE plot, i.e., the deviation of the instance-level effects from the average effect
- get a robust approximation of (a) the main ALE effect and (b) the heterogeneity, using automatic bin-splitting

References

- Paper repo: [git@github.com:givasile/RHALE.git](https://github.com:givasile/RHALE.git)
- Personal site: givasile.github.io
- Twitter: twitter.com/givasile1

