

Paper presentation at ACML 2022

DALE: Differential Accumulated Local Effects for efficient and accurate global explanations

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

¹ATHENA Research and Innovation Center

²Harokopio University of Athens

December 2022

Feature Effect

$y = f(x_s) \rightarrow$ plot showing the effect of x_s on the output y

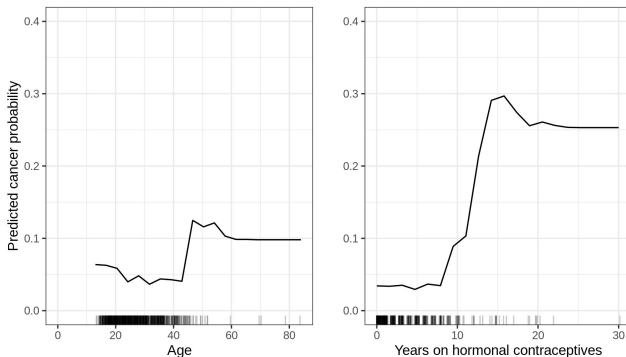


Figure: Image taken from Interpretable ML book.

Feature Effect is simple and intuitive.

Feature Effect Methods

- $x_s \rightarrow$ feature of interest, $\mathbf{x}_c \rightarrow$ other features
- FE methods take (f, \mathcal{D}, s) and return $y = f_{\langle \text{name} \rangle}(x_s)$
- PDP
 - ▶ Expected outcome over \mathbf{x}_c : $f(x_s) = \mathbb{E}_{\mathbf{x}_c}[f(x_s, \mathbf{x}_c)]$
 - ▶ **Unrealistic instances**
- MPlot
 - ▶ Expected outcome over $\mathbf{x}_c | x_s$: $f(x_s) = \mathbb{E}_{\mathbf{x}_c | x_s}[f(x_s, \mathbf{x}_c)]$
 - ▶ **Aggregated effects**
- ALE
 - ▶ $f(x_s) = \int_{x_{\min}}^{x_s} \mathbb{E}_{\mathbf{x}_c | z}[\frac{\partial f}{\partial x_s}(z, \mathbf{x}_c)] \partial z$
 - ▶ **Resolves both failure modes**

PDP vs MPlot vs ALE

DALE - Differential ALE

DALE, from the dataset $\mathcal{D} = \{\mathbf{x}^i, y^i\}_{i=1}^N$

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

- only change point effect computation
- Fast \rightarrow use of auto-differentiation, all derivatives in a single pass
- Versatile \rightarrow point effects computed once, change bins without cost
- Secure \rightarrow does not create artificial instances

For **differentiable** models, DALE resolves ALE weaknesses

DALE is faster and versatile

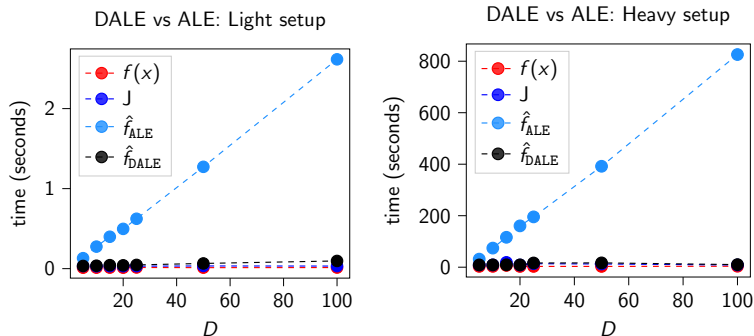
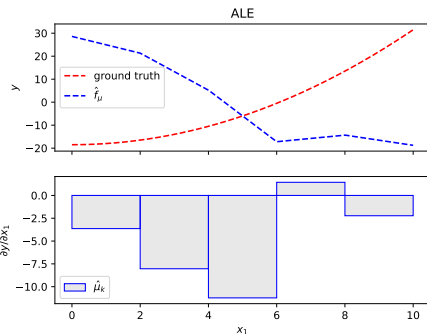
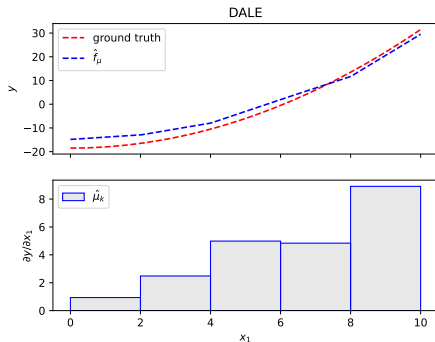


Figure: Light setup; small dataset ($N = 10^2$ instances), light f . Heavy setup; big dataset ($N = 10^5$ instances), heavy f

DALE considerably accelerates the estimation

DALE vs ALE - 5 Bins



- DALE: on-distribution, robust bin effect → good estimation
- ALE: completely OOD, robust bin effect → poor estimation