

# Bayesian Inference and Explainable AI

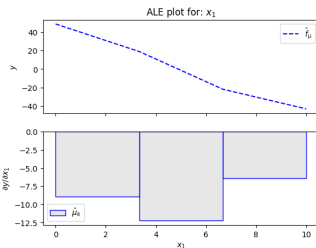
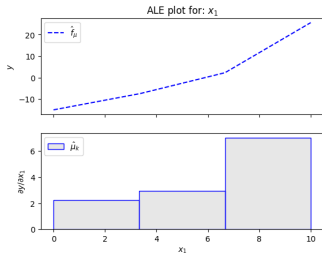
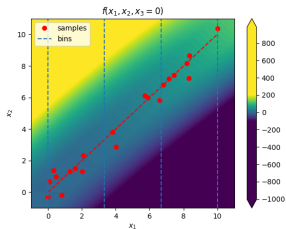
## Quantifying the uncertainty of the explanations

Vasilis Gkolemis<sup>1</sup>

<sup>1</sup>ATHENA Research and Innovation Center

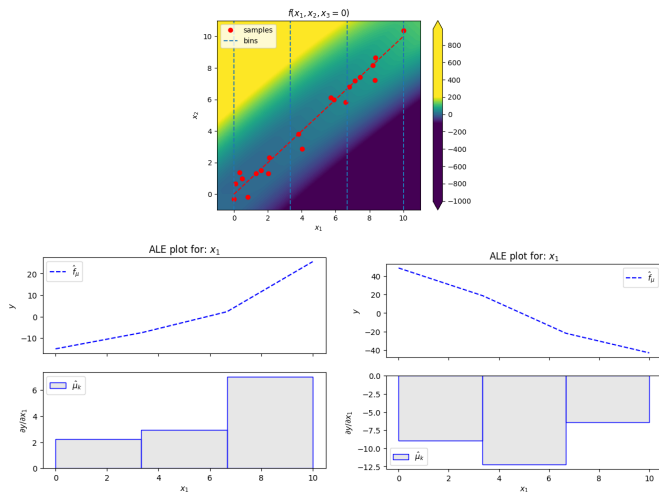
November 2021

# Feature Effect



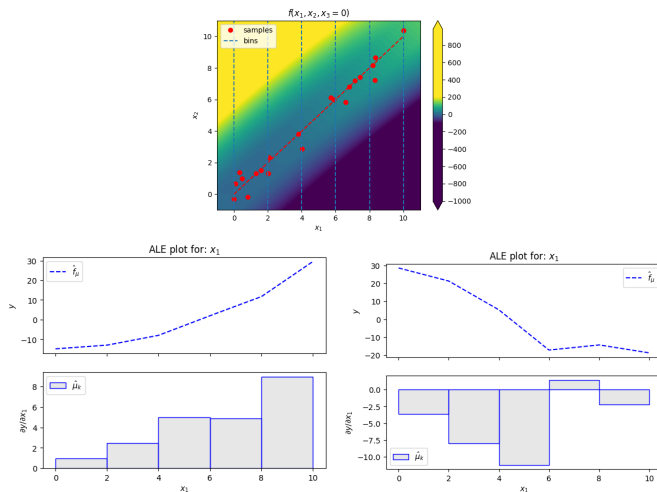
We search for a **single configuration (point-estimate)**  $\hat{\theta}$

# DALE vs ALE - 3 Bins



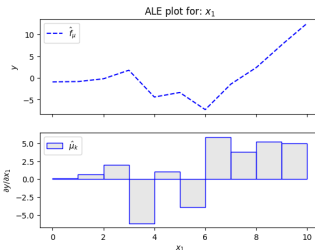
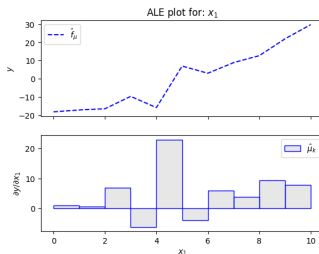
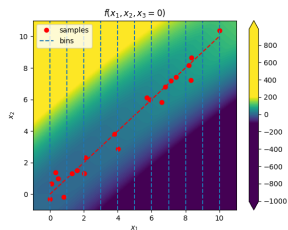
We search for a **single configuration (point-estimate)**  $\hat{\theta}$

## DALE vs ALE - 5 Bins



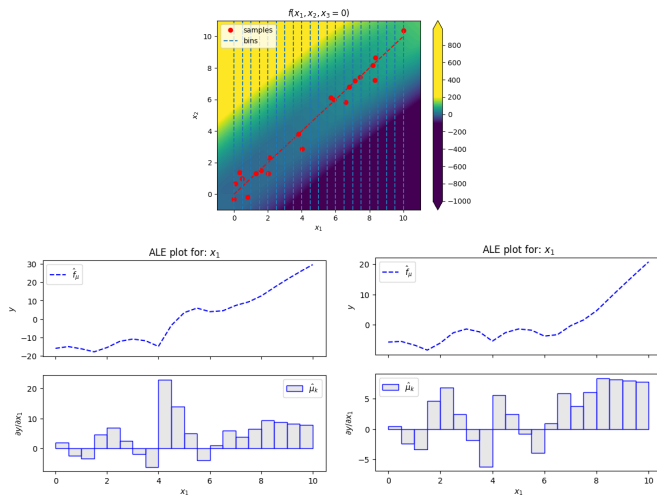
We search for a **single configuration (point-estimate)**  $\hat{\theta}$

## DALE vs ALE - 10 Bins



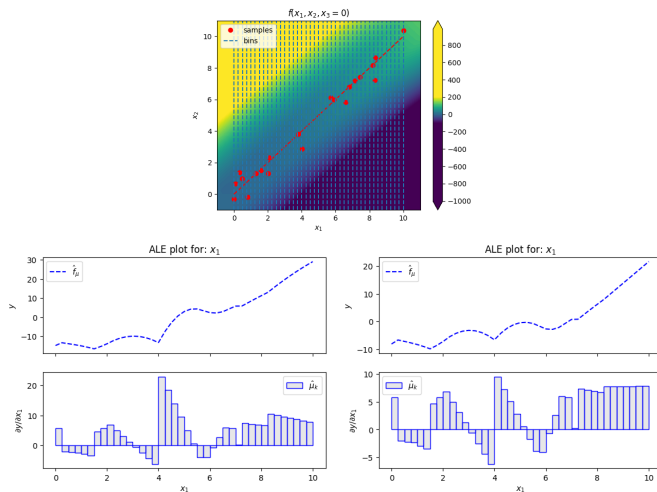
We search for a **single configuration (point-estimate)**  $\hat{\theta}$

## DALE vs ALE - 20 Bins



We search for a **single configuration (point-estimate)**  $\hat{\theta}$

## DALE vs ALE - 40 Bins



We search for a **single configuration (point-estimate)**  $\hat{\theta}$

# Traditional Machine Learning

- Parametric model  $f_{\theta} : \mathbf{x} \rightarrow y$
- Define a distance function  $d(\cdot, \cdot)$  and measure the distance (loss) from observed data

$$L(\theta) = \sum_i^N d(f_{\theta}(\mathbf{x}^i), y^i) \quad (1)$$

- Search for the parameter set  $\hat{\theta}$  that reproduces the observed data best

$$\hat{\theta} = \arg \min_{\theta} L(\theta) \quad (2)$$

We search for a **single configuration (point-estimate)**  $\hat{\theta}$



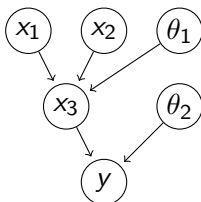
# Bayesian Formulation

- On the modelling part:
  - we need the joint distribution  $p(\mathbf{x}, y, \theta)$
  - to replace the parametric model  $f_\theta : \mathbf{x} \rightarrow y$
- Training part:
  - infer the posterior distribution  $p(\theta|D)$
  - to replace the optimal point estimate  $\hat{\theta} = \arg \min_{\theta} L(\theta)$
- Prediction part:
  - infer the predictive distribution  $p(y|\mathbf{x}, D)$
  - to replace the point-estimate prediction  $y = f_{\hat{\theta}}(\mathbf{x})$

We replace point estimates with distributions (uncertainty quantification)

# Modelling part

- joint distribution  $p(\mathbf{x}, y, \boldsymbol{\theta}) = p(y|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x})p(\boldsymbol{\theta})$
- $p(\boldsymbol{\theta})$ , our prior belief about the parameters of the model
- $p(y|\mathbf{x}, \boldsymbol{\theta})$ , the likelihood of the model
- joint distribution can be defined as a **DAG**



- We need to model  $p(\boldsymbol{\theta})$  and  $p(y|\mathbf{x}, \boldsymbol{\theta})$

# Training part

- We use Bayes law to infer the posterior distribution

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} \propto \prod_i^N p(y^i|\mathbf{x}^i, \theta)p(\theta) \quad (3)$$

where  $D = \{\mathbf{x}^i, y^i\}_{i=\{1, \dots, N\}}$ , the observed data (training-set)

- In the extreme case where  $p(\theta|D) = \delta(\theta - \hat{\theta})$ , we get a point-estimate is in traditional ML

The 'training process' leads to many possible models, each one with different probability (**uncertainty about the model**)

# Inference part

- We need to solve/approximate the predictive distribution  $p(y|\mathbf{x}, D) = \int_{\theta} p(y|\mathbf{x}, \theta) p(\theta|D) d\theta$
- We consider the posterior  $p(\theta|D)$  as known (computed exactly or approximated)
- In the extreme case where  $p(\theta|D) = \delta(\theta - \hat{\theta})$ , we get all the mass of the prediction  $p(y|\mathbf{x}, D) = p(y|\mathbf{x}, \hat{\theta})$  from a single model

The 'prediction process' gets one prediction per each plausible model (**uncertainty about the model leads to uncertainty about the prediction**)

# Bayesian Formulation - Disadvantages

What we lose

- On the modelling part
  - Time to think how the input features  $x_i$  relate to each other i.e. building the DAG
- On the training-prediction (inference) part
  - Expressions difficult to approximate
  - $p(\theta|D)$  - how to compute the posterior distribution?
  - $p(y|\mathbf{x}, D) = \int_{\theta} p(y|\mathbf{x}, \theta)p(\theta|D)d\theta$  - how to compute the predictive distribution?

Bayesian Formulation is **difficult from both the mathematical and the computational point-of-view**

# Bayesian Formulation - Advantages

What we get:

- On the modelling part
  - Specify who the features relate to each other
  - check [Model-based Machine Learning](#) - a new approach for ML model building
- On the inference part
  - Uncertainty estimation! Why we need it?
  - Most times the available data is not enough to reveal a single instance  $\hat{\theta}$
  - Sometimes we want to predict on a new  $x$  that is very different from the training set

Let's be wise enough and be uncertain about our predictions.

# Bayesian Formulation - Example

- In areas without training points our uncertainty is bigger