

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



# Scientific Inference with Interpretable Machine Learning

## Analyzing Models to Learn About Real-World Phenomena

Timo Freiesleben

Joint work with G. König, C. Molnar, & A. Tejero-Cantero

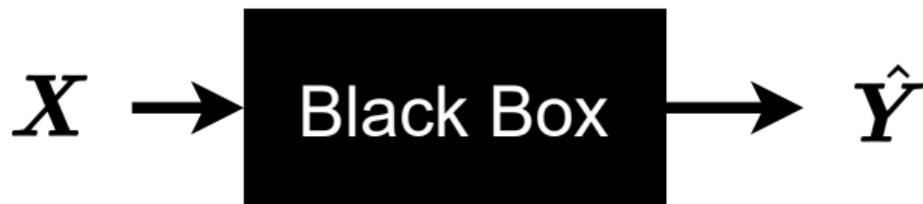
University of Tübingen, Cluster of Excellence 'Machine Learning for Science'

# Outline

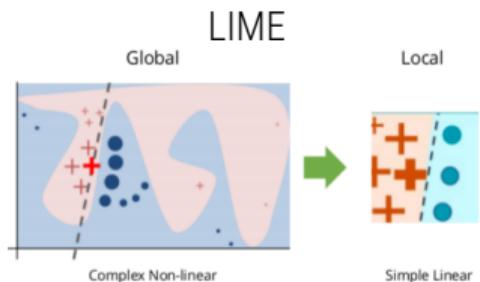
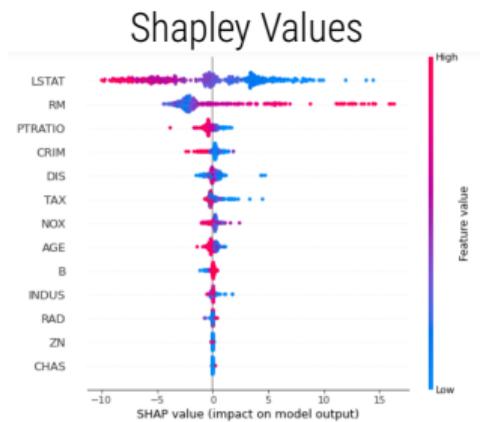
- 1 Motivation
- 2 Traditional scientific inference
- 3 Theory of property descriptors
- 4 Discussion

# Motivation: Interpretable ML

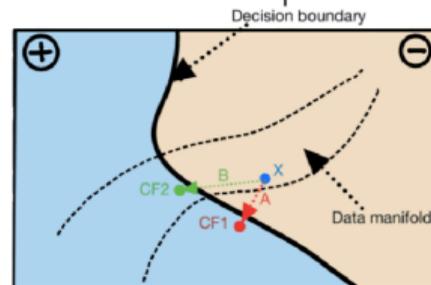
- ▶ Also called XAI
- ▶ Ingredients:
  - Data  $\mathcal{D}$
  - inputs  $\mathbf{X}$  and prediction  $\hat{\mathbf{Y}}$
  - Trained ML model  $\hat{m}$



# Motivation: Method zoo



## Counterfactual explanations

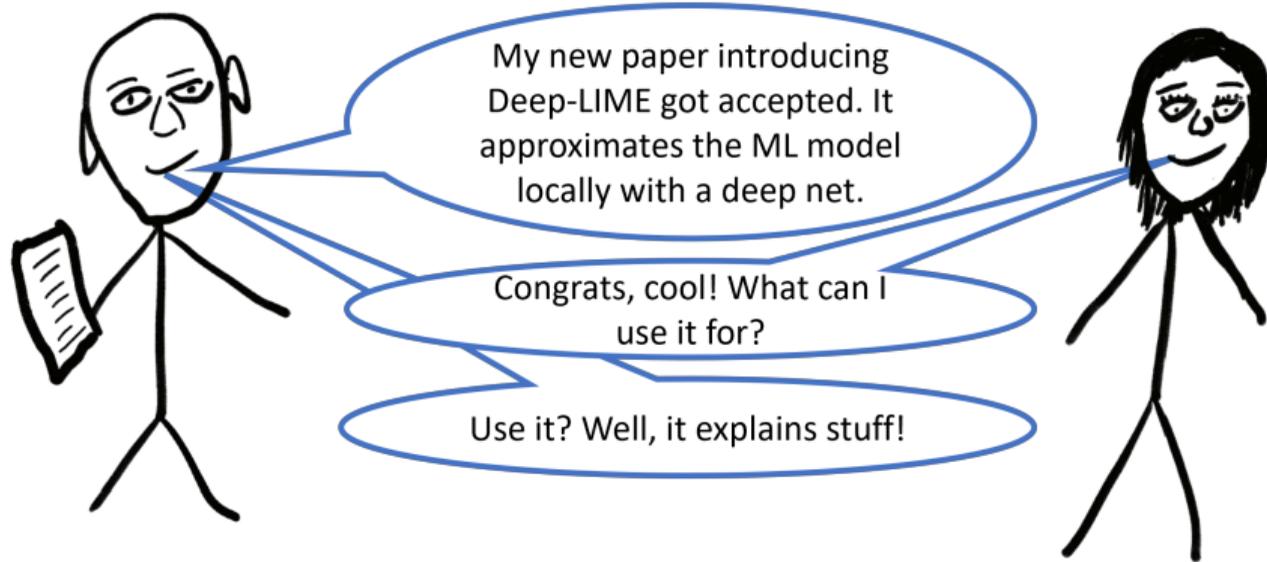


## Saliency maps



Images by: Idit Cohen, Mokuwe et al. [2020], Ribeiro et al. [2016], Verma et al. [2020]

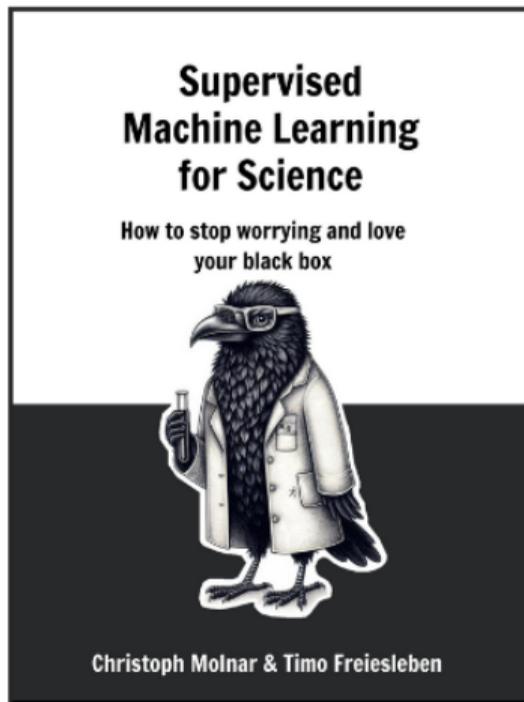
## Motivation: What real problems are solved?



Dear XAI community, we need to talk! [Freiesleben and König, 2023]

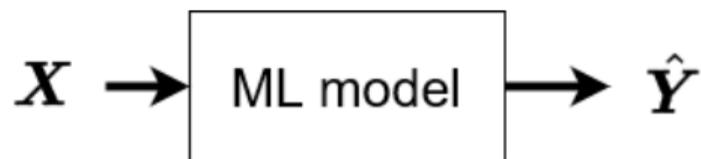
# Motivation: Can we use (interpretable) ML for science?

<https://ml-science-book.com/>

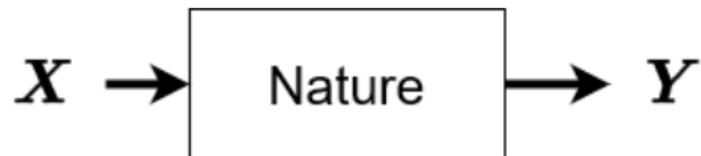


# Motivation: Model audit vs scientific inference

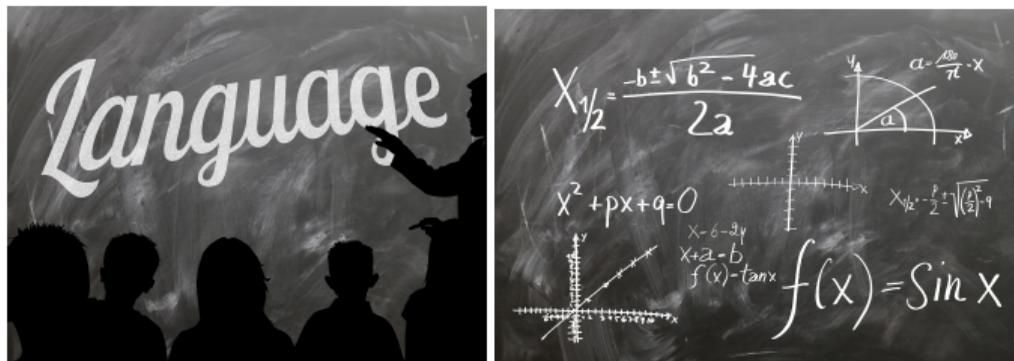
- ▶ Model audit



- ▶ Scientific inference



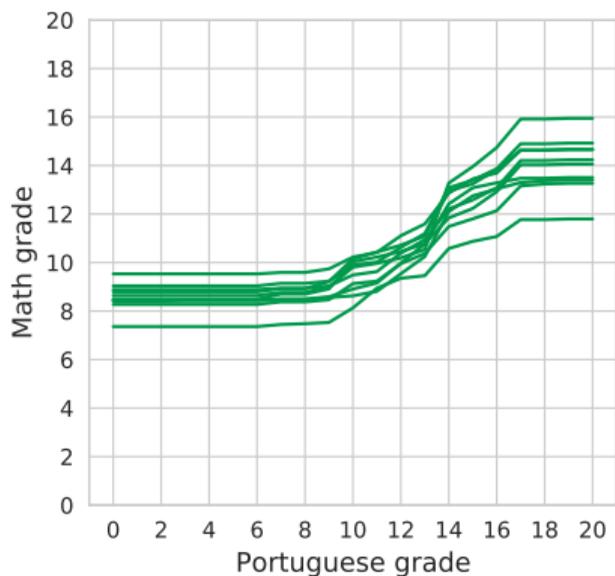
# Motivation: Laura example



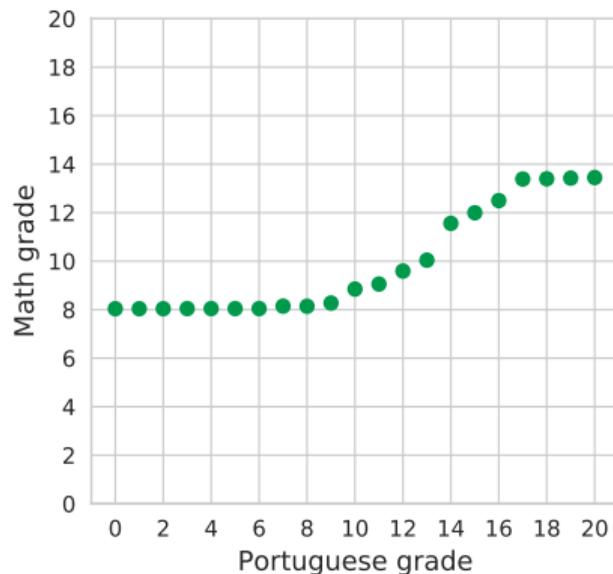
- ▶ How relate language and math skills?
- ▶ Data [Cortez and Silva, 2008]:
  - students grades,
  - parent's jobs/education,
  - age, tutoring, absences, etc.

# Motivation: Partial dependence plot

► Timo= (Port grade: 0; tutoring: *no*; absences: 5; ...)



ICE curves



Partial Dependence Plot (PDP)

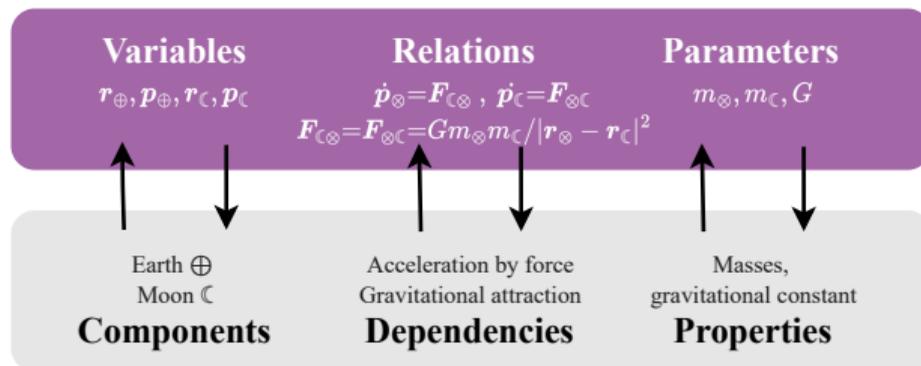
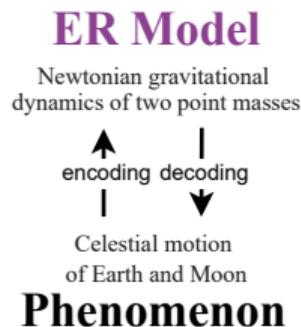
# Where Are We?

- 1 Motivation
- 2 Traditional scientific inference
- 3 Theory of property descriptors
- 4 Discussion

# Traditional scientific inference: Elementwise representation

## Definition: Elementwise Representationality

A model is *elementwise representational* (ER) if all model elements represent an element in the phenomenon.

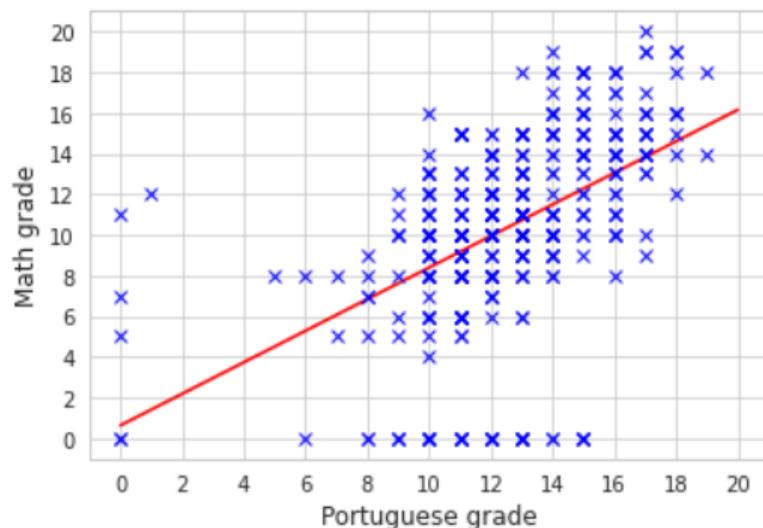


# Traditional scientific inference: Why ER?

- ▶ ER is cognitively appealing
- ▶ ER eases model construction
- ▶ **ER allows inference from model to world**

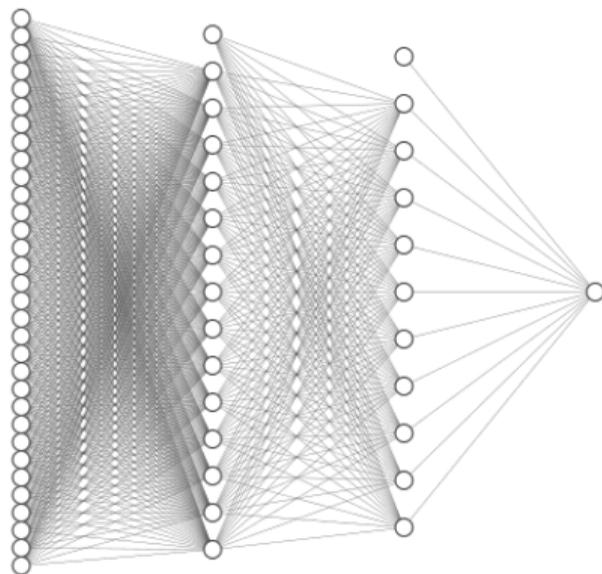
# Traditional scientific inference: Example

- ▶  $\mathbf{Y}$ =math and  $\mathbf{X}_p$ =Portuguese
- ▶  $\mathbf{Y} = \beta_0 + \beta_1\mathbf{X}_p + \epsilon$
- ▶ Least squares:  $\hat{m}_{LIN}(x_p) = 10.46 + 0.77x_p$
- ▶ Confidence Intervals  $CI_{\hat{\beta}_0} = [10.05; 10.88]$  and  $CI_{\hat{\beta}_1} = [0.63; 0.91]$



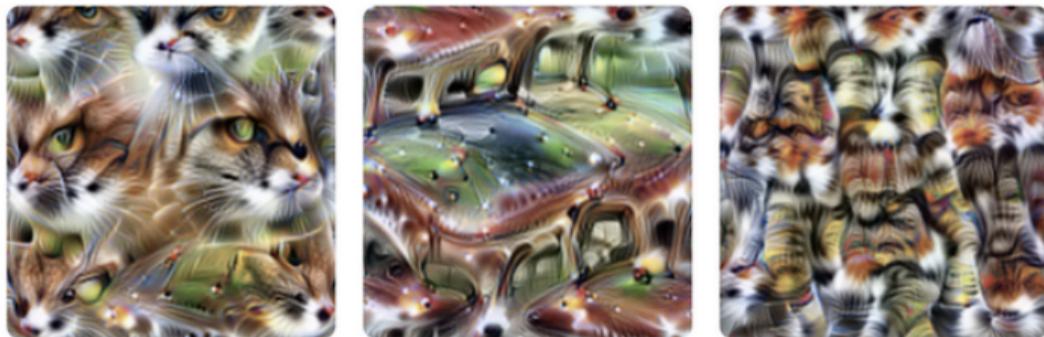
# Traditional scientific inference: ML models not ER

- ▶ ML models are less assumption laden
- ▶ Most model elements (weights, activation functions, etc) have no meaning



# Traditional scientific inference: Inference with ML

- ▶ Option 1: [Bokulich, 2011]
  - Scientific inference without ER is impossible
- ▶ Option 2: [Olah et al., 2020]
  - ML models are ER too

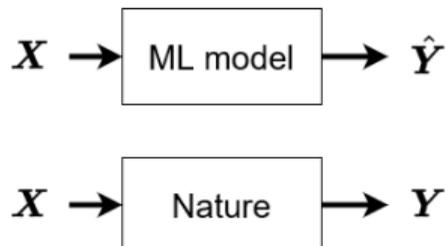


- ▶ Option 3: [Cichy and Kaiser, 2019]
  - IML for scientific inference

# Traditional scientific inference: IML for inference

## ► Problems

- Current IML focuses on model audit
- Not every audit allows for inference
- Audit and inference are complementary goals



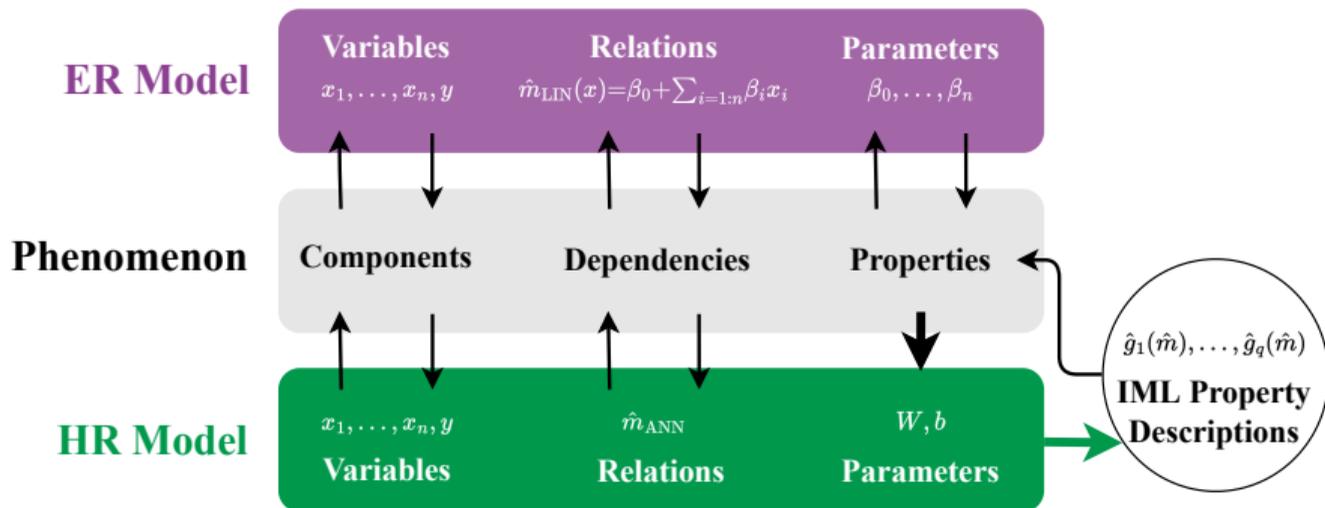
# Where Are We?

- 1 Motivation
- 2 Traditional scientific inference
- 3 Theory of property descriptors
- 4 Discussion

# Theory of property descriptors: Holistic representation

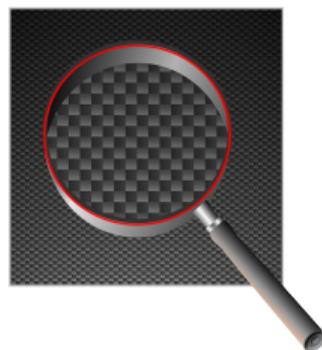
Definition: Holistic representationality

A model is *holistically representational* (HR) if the whole model represents aspects of the phenomenon.

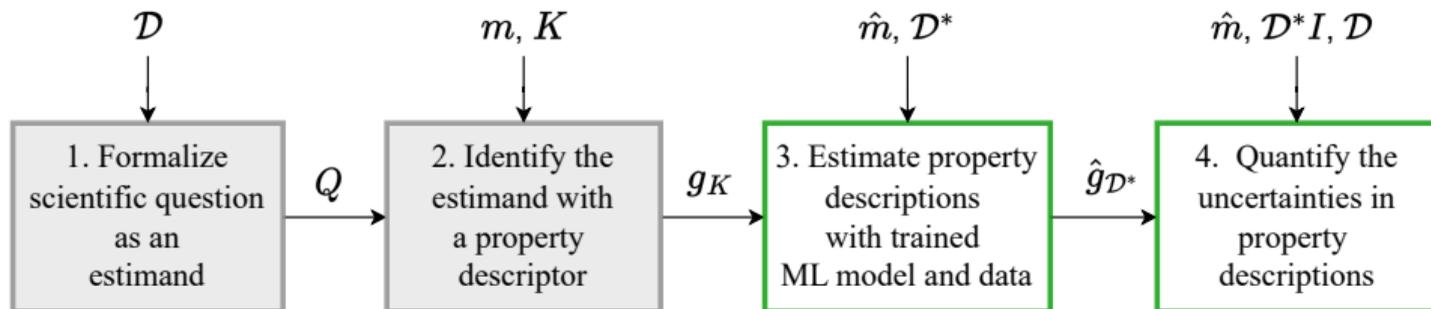


# Theory of property descriptors: What ML models represent?

<b>Problem</b>	<b>Loss</b>	$L(Y, \hat{m}(\mathbf{X}))$	<b>Optimal predictor<sup>a</sup><math>m</math></b>
Regression ( $Y$ continuous)	mean squared error	$(Y - \hat{m}(\mathbf{X}))^2$	$\mathbb{E}_{Y \mathbf{X}}[Y \mathbf{X}]$
	mean absolute error	$ Y - \hat{m}(\mathbf{X}) $	$\text{median}(Y \mathbf{X})$
Classification ( $Y$ discrete)	0-1 loss	0 if $\hat{m}(\mathbf{X}) = Y$ , else 1	$\arg \max_{y \in \mathcal{Y}} \mathbb{P}(Y=y \mathbf{X})$
	cross entropy	$\sum_{r \in \mathcal{Y}} \mathbb{P}_Y(r) \log \mathbb{P}_{\hat{m}(\mathbf{X})}(r)$	$\mathbb{P}(Y \mathbf{X})$



# Theory of property descriptors: Four steps



# Theory of property descriptors: 1. Formalize scientific question

- ▶ Scientists start by asking and formalizing questions.
- ▶ Question: How are language skills associated with math skills?
- ▶ Formalized Question:  $Q = \mathbb{E}_{\mathbf{Y}|\mathbf{X}_p}[\mathbf{Y} | \mathbf{X}_p]$

## Theory of property descriptors: 2. Identify estimand

### Definition: Question Identifiability

We say that a question is *identifiable* relative to probabilistic knowledge  $K$  if we can compute  $Q$  from  $m$  and  $K$ .

- ▶ Laura's question can be identified with  $K = \mathbb{P}(\mathbf{X}_{-p} \mid \mathbf{X}_p)$

$$\begin{aligned} Q &:= \mathbb{E}_{\mathbf{Y} \mid \mathbf{X}_p}[\mathbf{Y} \mid \mathbf{X}_p] \\ &= \mathbb{E}_{\mathbf{X}_{-p} \mid \mathbf{X}_p}[\mathbb{E}_{\mathbf{Y} \mid \mathbf{X}}[\mathbf{Y} \mid \mathbf{X}] \mid \mathbf{X}_p] && \text{(by the tower rule)} \\ &= \mathbb{E}_{\mathbf{X}_{-p} \mid \mathbf{X}_p}[m(\mathbf{X}) \mid \mathbf{X}_p]. \end{aligned}$$

## Theory of property descriptors: 2. Identify estimand

### Definition: Property Descriptor

A *property descriptor* is a continuous function  $g_K$  that identifies  $Q$  given  $K$

$$g_K : \mathcal{M} \rightarrow \mathcal{Q} \quad \text{with} \quad g_K(m) = Q.$$

- In our example, this is the conditional Partial Dependence Plot (cPDP):

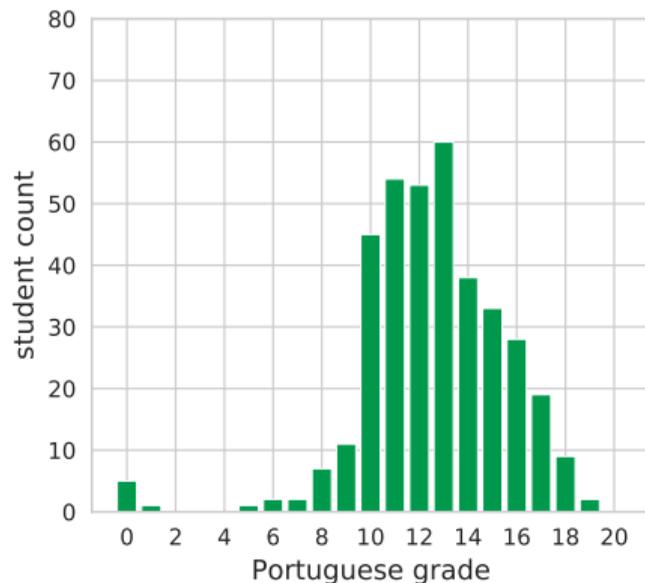
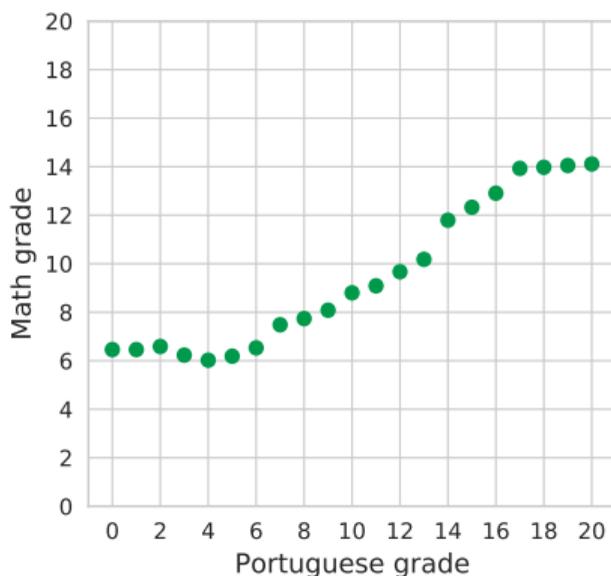
$$g_K(\hat{m}) := \mathbb{E}_{\mathbf{X}_{-p} | \mathbf{x}_p} [\hat{m}(\mathbf{X}) | \mathbf{X}_p]$$

# Theory of property descriptors: 3. Estimate property

In real life, we have limited access to  $\mathbf{X}$ ,  $\mathbf{Y}$ . We have finite data.

Definition: Property Description Estimator

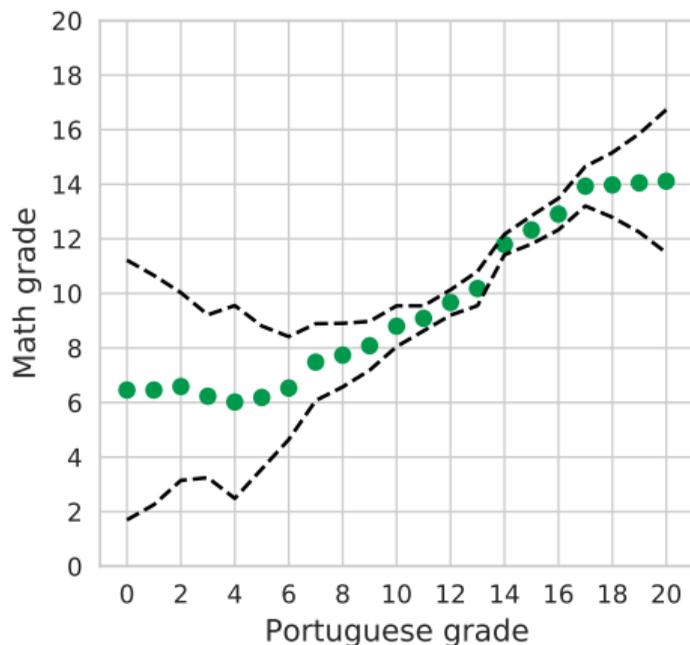
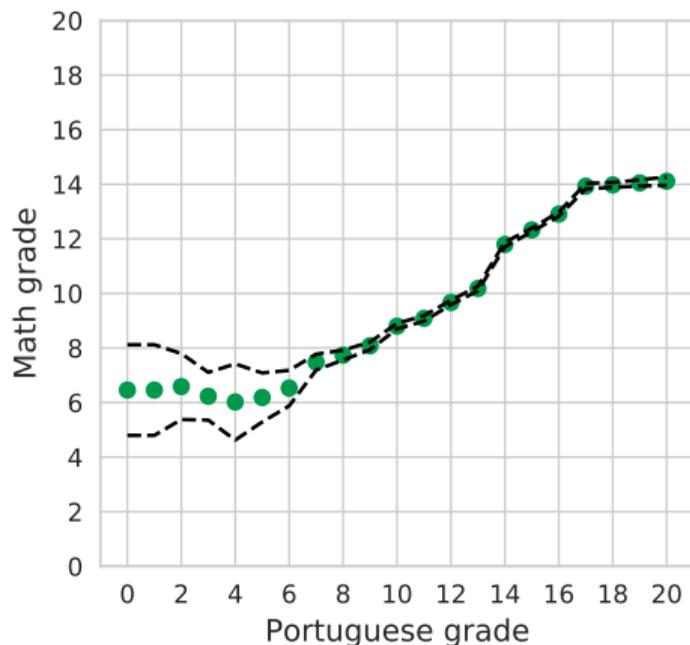
The *property description estimator*  $\hat{g}_{\mathcal{D}^*}$  is an unbiased estimator of  $g_K$ .



# Theory of property descriptors: 4. Uncertainty quantification

We make two errors on the way:

- 1 we do not have the optimal model (model error), and
- 2 we only have finite data (estimation error).



# Theory of property descriptors: Practical descriptors

	Global / local question	Estimand	IML method
Conditional contribution	How much worse can $Y$ be predicted from $X$ if we did not know $X_p$ ?	$EPE_{\mathbf{X},Y} m_{\mathbf{X}}(X) - EPE_{\mathbf{X}_{-p},Y} m_{\mathbf{X}_{-p}}(\mathbf{X}_{-p})$	cFI Strobl et al. (2008)
	How much worse can $Y$ be predicted from $X = \mathbf{x}$ if we did not know $X_p$ ?	$L(y, m_{\mathbf{X}}(\mathbf{x})) - L(y, m_{\mathbf{X}_{-p}}(\mathbf{x}_{-p}))$	ICI Casalicchio et al. (2019)
Fair contribution	What is the fair share of feature $X_p$ in the prediction of $Y$ ?	$\frac{1}{n} \sum_{S \subset \mathcal{N} \setminus \{p\}} \binom{n-1}{ S }^{-1} (EPE_{\mathbf{X}_{S \cup \{p\}}, Y} m_{\mathbf{X}_{S \cup \{p\}}}(\mathbf{X}_{S \cup \{p\}}) - EPE_{\mathbf{X}_S, Y} m_{\mathbf{X}_S}(\mathbf{X}_S))$	SAGE Covert et al. (2020)
	What is the fair share of feature $X_p$ in the prediction of $Y$ if $X = \mathbf{x}$ ?	$\frac{1}{n} \sum_{S \subset \mathcal{N} \setminus \{p\}} \binom{n-1}{ S }^{-1} (m_{\mathbf{X}_{S \cup \{p\}}}(\mathbf{x}_S, x_p) - m_{\mathbf{X}_S}(\mathbf{x}_S))$	Conditional Shapley values Aas et al. (2021)

- Which student information should educators track?  $\Rightarrow$  ICI, cFI, conditional SHAP & SAGE

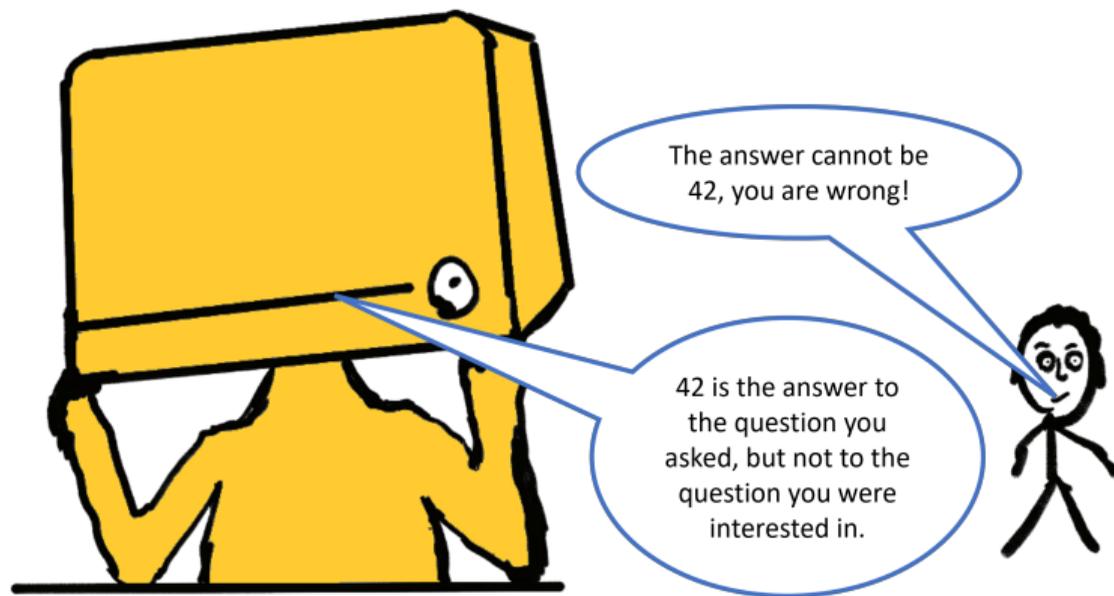
# Theory of property descriptors: Practical descriptors

	Global / local question	Estimand	IML method
Effect	What is the best estimate of $Y$ if we only know $X_p$ ?	$m_{X_p}(X_p)$	cPDP <small>Apley and Zhu (2020)</small>
	How does the best estimate of $Y$ change relative to $X_p$ , knowing that $X_{-p} = x_{-p}$ ?	$m_{\mathcal{X}}(X_p, x_{-p})$	ICE curve <small>Goldstein et al. (2015)</small>
Relevant value	Under which realistic conditions $\mathcal{X}$ can we observe relevant value $y_{rel}$ ?	$\arg \min_{x \in \text{supp } \mathcal{X}} d_Y(m_{\mathcal{X}}(x), y_{rel})$	PRIM <sup>a</sup> <small>Friedman and Fisher (1999)</small>
	Under which realistic conditions similar to $x$ can we observe relevant value $y_{rel}$ ?	$\arg \min_{x' \in \text{supp } \mathcal{X}} d_Y(m_{\mathcal{X}}(x'), y_{rel}) + \lambda d_{\mathcal{X}}(x, x')$	Counterfactuals <sup>b</sup> <small>Dandl et al. (2020)</small>

- ▶ How influences parents' education students math skills?  $\Rightarrow$  ICE & cPDP
- ▶ What characterizes (more/less) successful students?  $\Rightarrow$  counterfactuals & PRIM

# Theory of property descriptors: Disagreement

- ▶ Methods can only meaningfully disagree if they have different estimands.
- ▶ The disagreement problem stems from a lack of clarity about the question asked.



# Where Are We?

- 1 Motivation
- 2 Traditional scientific inference
- 3 Theory of property descriptors
- 4 Discussion

# Discussion: Causality

- ▶ Most scientific questions are causal:
  - Would tutoring in Portuguese improve students math skills? (**Interventional**)
  - Did the student fail in math because of her Portuguese skills? (**Counterfactual**)
- ▶ Property descriptors describe associational quantities.
- ▶ Causal questions add another layer to the pipeline and require causal knowledge.

## Discussion: Limitations

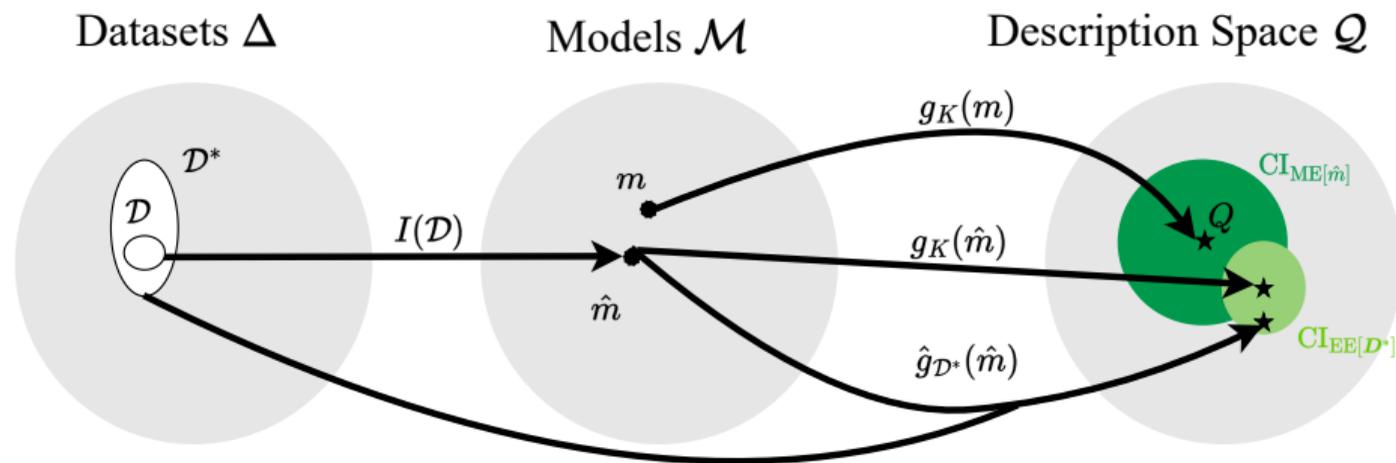
- ▶ Direct estimation often better! (e.g. targeted learning [[Van der Laan and Rose, 2011](#)])
- ▶ Conditional sampling is needed but hard!
- ▶ Formalizing questions on images and sound?

- ▶ **Problem:** Scientific inference via model elements is not available. Current IML mixes different desiderata.
- ▶ **Our Solution:** Smart interrogation with property descriptors allows to learn about the process.

# Questions



# Graph for Formal Depiction



- A. Bokulich. How scientific models can explain. *Synthese*, 180(1):33–45, 2011.
- R. M. Cichy and D. Kaiser. Deep neural networks as scientific models. *Trends in cognitive sciences*, 23(4):305–317, 2019.
- P. Cortez and A. M. G. Silva. Using data mining to predict secondary school student performance. 2008.
- T. Freiesleben and G. König. Dear xai community, we need to talk! fundamental misconceptions in current xai research. In *World Conference on Explainable Artificial Intelligence*, pages 48–65. Springer, 2023.
- M. Mokuwe, M. Burke, and A. S. Bosman. Black-box saliency map generation using bayesian optimisation. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020.
- C. Olah, N. Cammarata, L. Schubert, G. Goh, M. Petrov, and S. Carter. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- M. T. Ribeiro, S. Singh, and C. Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016. doi: 10.1145/2939672.2939778.
- M. J. Van der Laan and S. Rose. *Targeted learning*. Springer, 2011.

S. Verma, J. Dickerson, and K. Hines. Counterfactual explanations for machine learning: A review. *arXiv preprint arXiv:2010.10596*, 2, 2020.