Maria Moreno de Castro
moreno@dkrz.de

# Interpretable Machine Learning

DKRZ

# ML algorithms encrypt and magnify bias

## IKEA Effect

We place higher value on things we partially created ourselves.

## Blind Spot Bias

We don't think we have bias, and we see it in others more than ourselves.

*"I am not biased!"*

## Automation Bias

We rely on automated systems, sometimes trusting too much in the automated correction of actually correct decisions.

## Law of Triviality (aka "Bike-Shedding")

We give disproportionate weight to trivial issues, often while avoiding more complex issues.

in geosciences…

- large spread in climate prediction is often nowadays attributed to cloud radiative effects → resolving clouds seems to be the way to go to reduce epistemic uncertainty, however the ecosystem modelling community often expresses that the GCM are biased because they do not resolve phytoplankton and in general underrepresent the biological component

- geographical areas we do not have as much observational data, modelling groups,...

- use ensemble means as it everything would be normally distributed

- use deterministic solvers for differential equations with stochastic components

- our physical models are built in blocks by different people at different institutions at different times, with tons of simplifications, approximations, assumptions, and empirical parameters and to make all work together we use tuning parameters which do not have physical meaning and introduce compensation errors

ML algorithms can easily break the

# the fundamental laws of physics

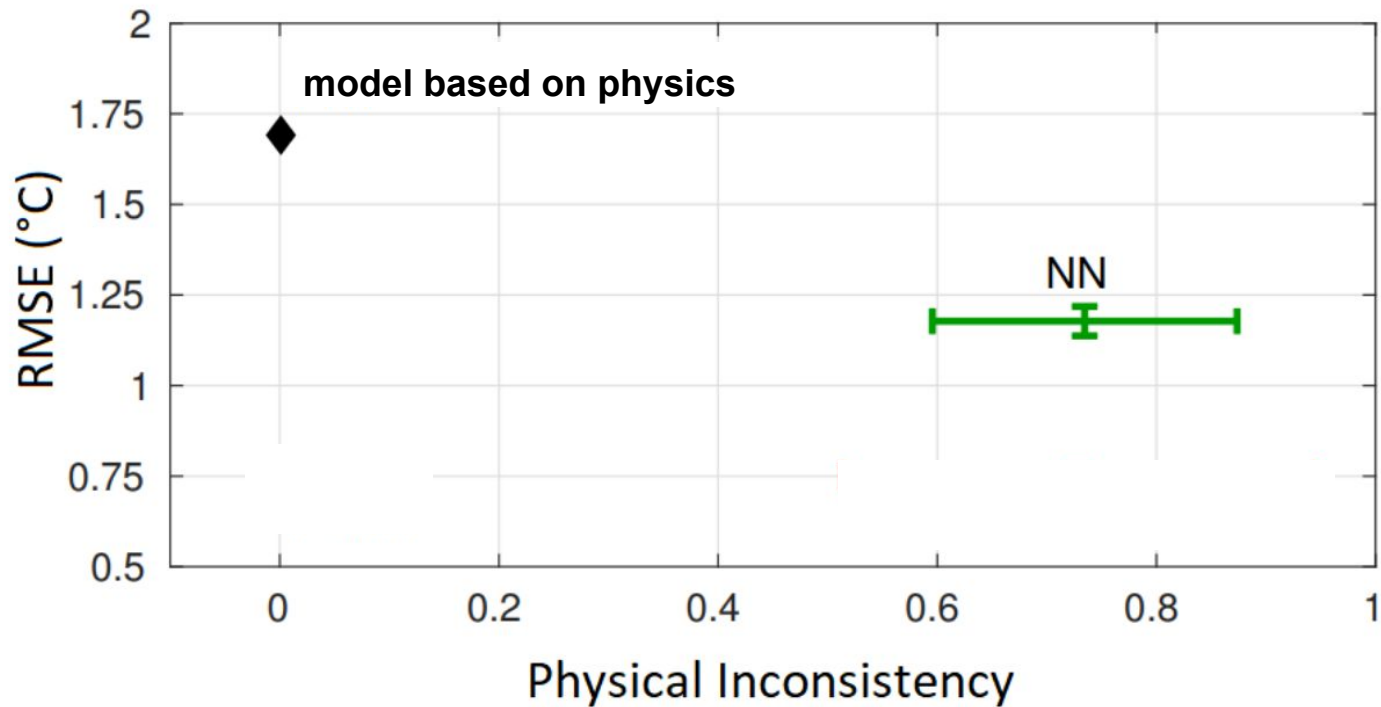energy and mass conservations, nonnegative densities, precipitations,...



[Noether's theorem](#) explains why [conservation laws](#) exists (wikipedia)
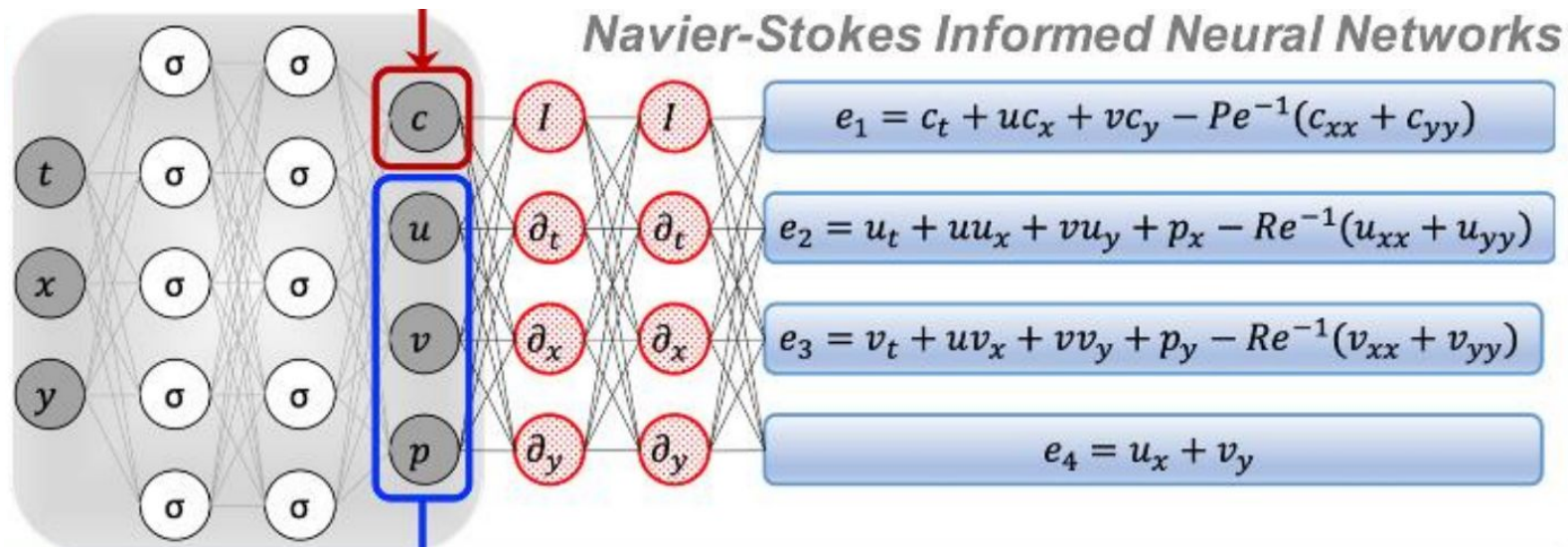
Example:
lakes simulations

| Depth (m) | Temp (°C) |
|-----------|-----------|
|           |           |
|           |           |



Results on Lake Mille Lacs

Example:
fluid dynamics



Navier-Stokes Informed Neural Networks

$$e_1 = c_t + uc_x + vc_y - Pe^{-1}(c_{xx} + c_{yy})$$

$$e_2 = u_t + uu_x + vu_y + p_x - Re^{-1}(u_{xx} + u_{yy})$$

$$e_3 = v_t + uv_x + vv_y + p_y - Re^{-1}(v_{xx} + v_{yy})$$

$$e_4 = u_x + v_y$$

- **Black-box models**

Humans cannot understand the cause of the decisions: knowing the value of the parameters is not enough to infer what is going on and/or underlying assumptions/limitations are unknown so it is hard to spot when these models are biased.

Examples: Random forest, NN,...

- **Black-box models**

Humans cannot understand the cause of the decisions: knowing the value of the parameters is not enough to infer what is going on and/or underlying assumptions/limitations are unknown so it is hard to spot when these models are biased.

Examples: Random forest, NN,...

- **Interpretable models or Glass-box models**

Humans can understand the cause of a decision: knowing the value of the parameters helps and the underlying assumptions/limitations are known.

Examples: linear models, logistic regression, decision trees, naive Bayes, and k-nearest neighbors.

- **Explainable models**

The models are still black-boxes but we use some methods (based on surrogate models) a posteriori to try to infer where the predictions came from.
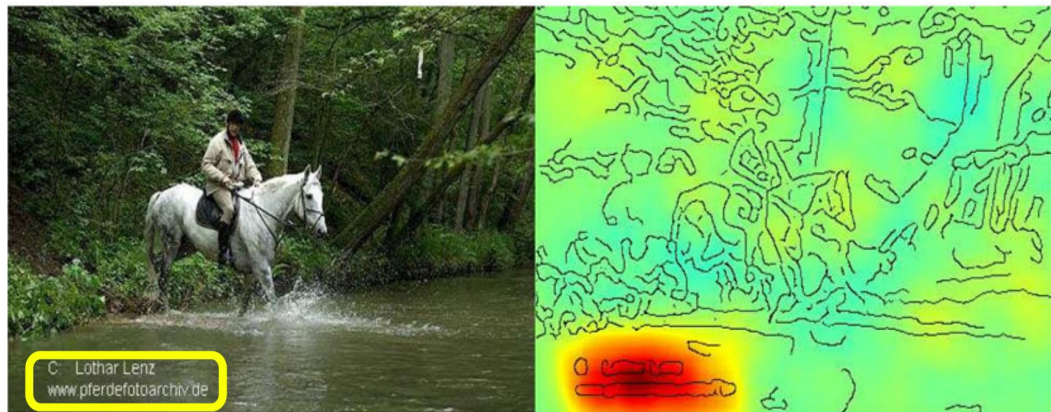
- **Explainable models**

The models are still black-boxes but we use some methods (based on surrogate models) a posteriori to try to infer where the predictions came from.

- sensitivity analysis based on observing the effect of perturbations on model components

- identify what features or feature values contributed the most to the predictions (often presented in saliency maps)

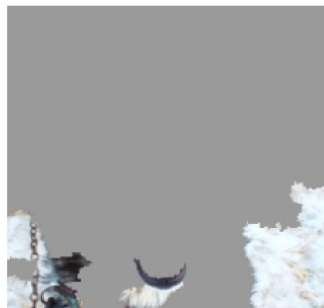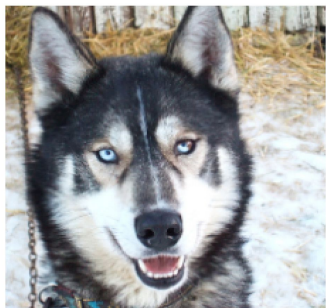- for every model, many available libraries (skitlearn, Tensorflow,... or extensions)

✓ essential        ✓ feasible        ✓ add scientific value

# Layer-wise
# Relevance Propagation (LRP)



single prediction, run again a back propagation that **tracks** what was activated and how much (nice but **costly**)

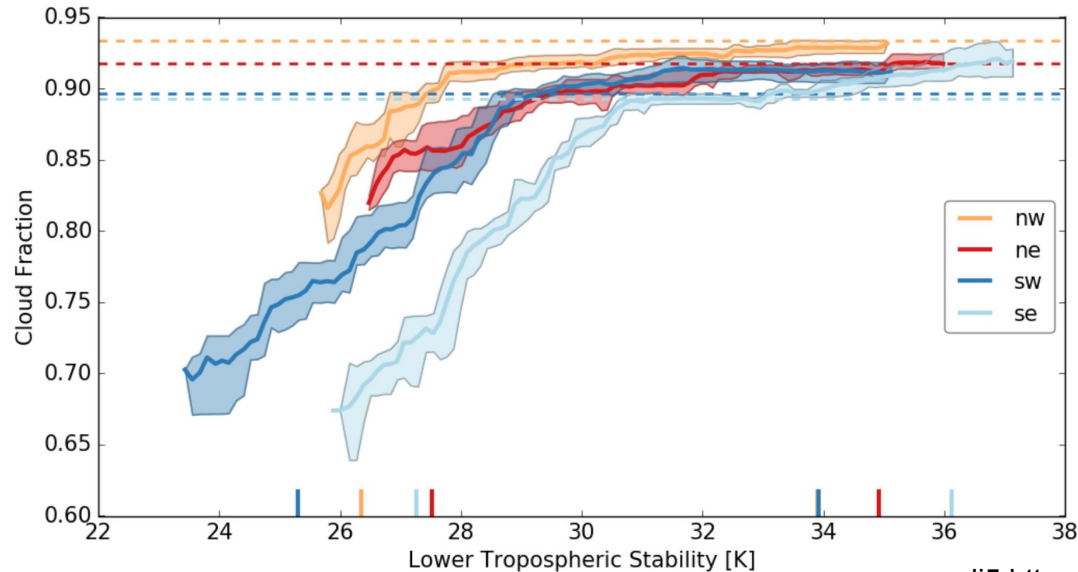# Local Interpretable Model-agnostic Explanation (LIME)



single prediction, run an interpretable model with **the black-box prediction as target** (**superpixels**)
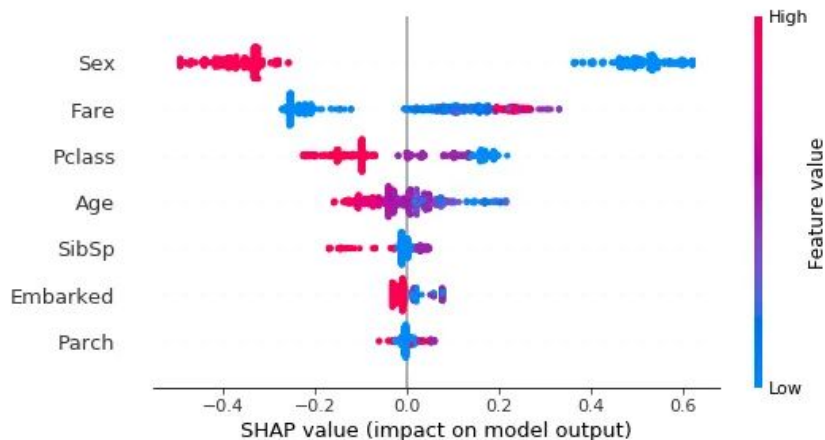
# Feature importance

**Permutation importance**: for any supervised model, just **shuffle** a feature of the validation data and compare performance (what **features values** are important?)

**Partial dependence plots**: 1. select feature 2. define a grid 3. per grid value: a) **replace feature with grid value** and b) average prediction 4. draw curve. (**cancelation effects? how to know the direction of the effect?**)



Fuchs et al. 2018 (ACP)

eli5 https://eli5.readthedocs.io/en/latest/,
PDPBox https://pdpbox.readthedocs.io/en/latest/,...

# Shapley Additive exPlanation (SHAP)

```
# summarize the effects of all the features
shap.summary_plot(shap_values, X)
```



1. single prediction, based on game theory, a prediction can be explained by assuming that each feature value of the instance is a "player" in a game where the prediction of that instance is the "payout", in a linear world is easy, just adding the contributions, but here **coalitions, synergies, direction of the effect**,... are considered.

2. do it for all of the prediction to get the global shapley values (**costly**)

https://github.com/slundberg/shap

# warning

XAI techniques are not the ultimate solution: the surrogate models bring their own assumptions and limitations, and are error-prone, an interpretable model is always more trustable

## Fooling LIME and SHAP: Adversarial Attacks on Post hoc Explanation Methods

**Authors:** Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, Himabindu Lakkaraju

Authors Info & Affiliations

# Parameterization schemes

# Interpreting and Stabilizing Machine-learning Parametrizations of Convection

Noah D. Brenowitz[1], Tom Beucler[2,3], Michael Pritchard[2], and Christopher S. Bretherton[1]

Predicting the behavior of NNs is tied to the difficult problem of interpreting NN emulators of physical processes. [...] How can we tailor ML interpretability techniques [...] for the particular purpose of interpreting NN parameterizations of convection?

1) Partial dependence plots to tests the nonlinear sensitivity of a single ML parametrization to systematic changes in its inputs

2) Linear Response Functions (LRF) or saliency map extend the analysis to the full input space of a parameterization. In a previous paper they computed LRFs to analyze what was causing their NN parameterizations to produce unstable dynamics when coupled to a GCM → reducing the potential for spurious causality by ablating both the upper atmospheric temperature and humidity from the input features of an NN parameterization results in a stable scheme [...] which demonstrates that ML interpretability techniques have already significantly aided the development of ML parameterizations.

# "Dendrology" in Numerical Weather Prediction: What Random Forests and Logistic Regression Tell Us about Forecasting Extreme Precipitation

GREGORY R. HERMAN AND RUSS S. SCHUMACHER

Feature importance for Random Forest or Gini importance:

the number of splits based on the feature summed over the forest [...], it is normalized, an importance of one then indicates that all decision nodes in every tree of the forest split on the corresponding feature, while an importance of zero indicates that no decision node splits based on that feature.

Permutation accuracy importance:

for each predictive feature, the feature value for each sample used to construct a given tree is permuted to a different sample's value. Importance is determined by the decline in the model's predictive performance when replacing the true values with the permuted ones.

# Using Machine Learning to Parameterize Moist Convection: Potential for Modeling of Climate, Climate Change, and Extreme Events

Paul A. O'Gorman[1] (iD) and John G. Dwyer[1] (iD)

1. use the Random Forest parameterization to generate a Linear Response Function for the response of convective precipitation to small perturbations in temperature, specific humidity, and surface pressure.

2. use the concept of feature importance which seeks to measure the importance of the different input features (here temperature and humidity at different levels and surface pressure [...] for both the occurrence and strength of convection. [...] is implemented in RandomForestRegressor class of scikit-learn.

   [...] the feature importance is a useful additional diagnostic for the interaction of convection with the large-scale environment.

# Applications of Deep Learning to Ocean Data Inference and Subgrid Parameterization

Thomas Bolton[1] and Laure Zanna[1]

Machine Learning for Stochastic Parameterization:
Generative Adversarial Networks in the Lorenz '96 Model

# Could Machine Learning Break the Convection Parameterization Deadlock?

P. Gentine[1], M. Pritchard[2], S. Rasp[3], G. Reinaudi[1], and G. Yacalis[2]

# Deep learning to represent subgrid processes in climate models

Stephan Rasp[a,b,1], Michael S. Pritchard[b], and Pierre Gentine[c,d]

Published: 22 July 2017

Parameterization of typhoon-induced ocean cooling using temperature equation and machine learning algorithms: an example of typhoon Soulik (2013)

Jun Wei ✉, Guo-Qing Jiang & Xin Liu

Research Letter | 🔒 Open Access | cc ⓘ ⊜ Ⓢ

A Deep Learning Algorithm of Neural Network for the Parameterization of Typhoon-Ocean Feedback in Typhoon Forecast Models

Guo-Qing Jiang, Jing Xu, Jun Wei ✉

# References

- Bias
  - Google face recognition: https://twitter.com/jackyalcine/status/615329515909156865 and https://www.wired.com/story/when-it-comes-to-gorillas-google-photos-remains-blind/
  - Man in kitchen: https://jyzhao.net/ and https://searchenterpriseai.techtarget.com/feature/Big-data-throws-big-biases-into-machine-learning-data-sets
  - 50 types https://www.visualcapitalist.com/50-cognitive-biases-in-the-modern-world/
- Physics-guided neural networks
  - Lakes simulation: Karpatne et al 2018 https://arxiv.org/pdf/1710.11431.pdf and https://towardsdatascience.com/physics-guided-neural-networks-pgnns-8fe9dbad9414
  - Navier-Stokes: Raissi et al 2020 https://www.sciencedirect.com/science/article/pii/S0021999118307125
- Uncertainty
  - Epistemic: http://yingzhenli.net/home/pdf/epistemic_uncertainty_neurips_bdl2019.pdf
  - Under data shift: Ovadia et al 2020 https://papers.nips.cc/paper/9547-can-you-trust-your-models-uncertainty-evaluating-predictive-uncertainty-under-dataset-shift.pdf
- XAI
  - Horse and LRP: https://www.nature.com/articles/s41467-019-08987-4
  - Husky vs Wolf and LIME: https://arxiv.org/pdf/1602.04938.pdf
  - Feature importance, partial dependence plots, and individual conditional expectation https://www.kaggle.com/learn/machine-learning-explainability

# WEAPONS OF MATH DESTRUCTION

HOW BIG DATA INCREASES INEQUALITY
AND THREATENS DEMOCRACY

## CATHY O'NEIL

'Wise, fierce and desperately necessary'
JORDAN ELLENBERG

# Interpretable
# Machine Learning

### A Guide for Making
### Black Box Models Explainable

@ChristophMolnar
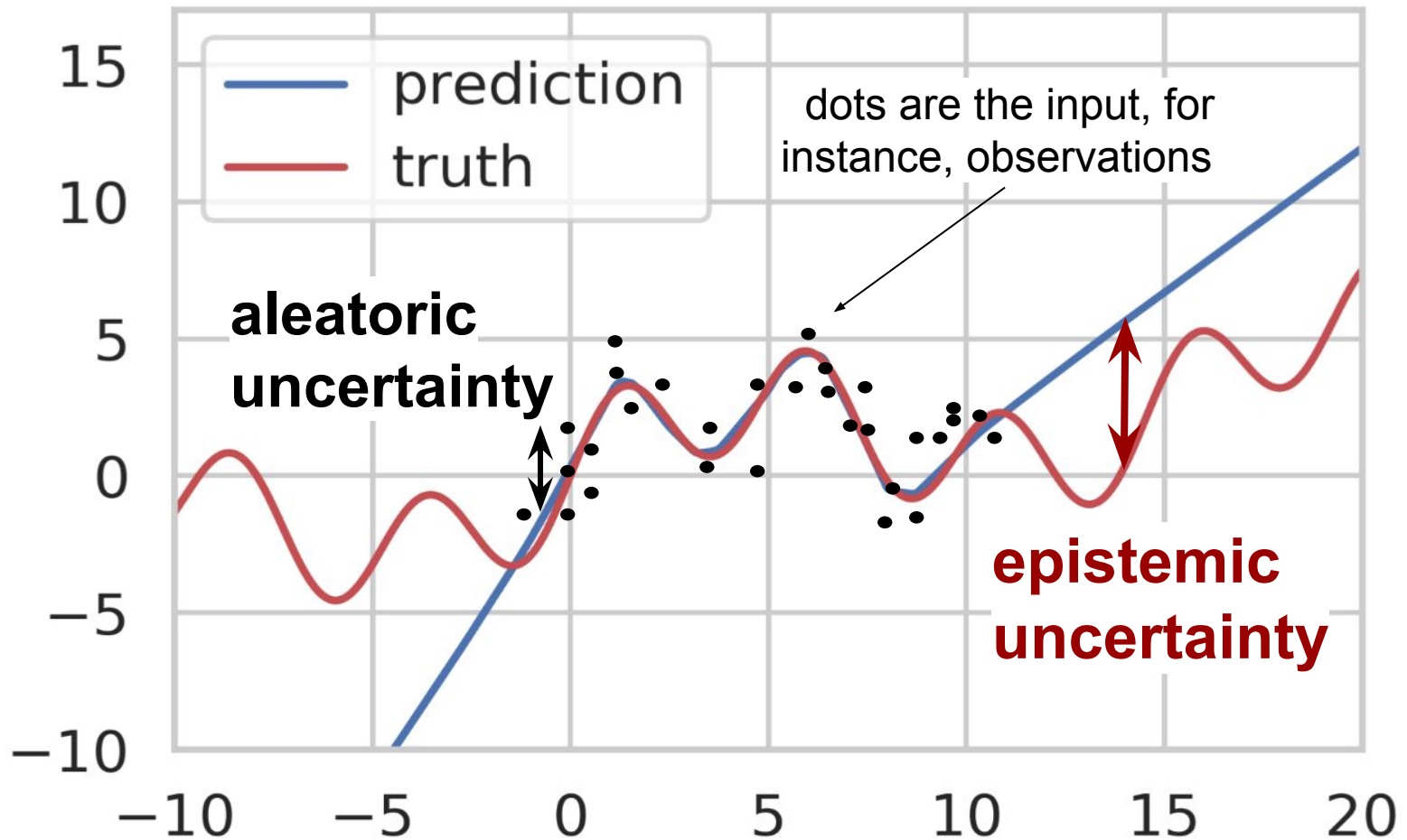
https://christophm.github.io
/interpretable-ml-book/

# Two types of uncertainty

Aleatoric: "what is the next outcome of tossing a coin?" related to an individual experiment outcome, it is non-reducible with more input data, it is the noise in the data.

Epistemic: "How much do I believe the coin is fair?" it is related to the model's belief after seeing the sample, it does reduce when having more data.

moreno@dkrz.de

ML are designed for interpolation, not extrapolation

**aleatoric uncertainty**

dots are the input, for instance, observations

**epistemic uncertainty**

# solutions:
## Gaussian Processes, Monte Carlo dropout, deep ensembles, dropout ensembles, and quantile regression

Florian Wilhelm: Are you sure about that?!
Uncertainty Quantification in AI | PyData...
PyData • 162 views • 1 month ago

Actually, there is a 3rd type of uncertainty:

Distribution shift: "Am I still flipping the same coin?" it is related to changes of the underlying quantity of interest, we assume that training and application data are i.i.d. but data drifts in time, we apply the model to data from a different location, the labeller changed,...

many problems in geoscience are non stationary

- training data are not longer representative if the system has changed
- the accuracy of the trained model definitely decreased under data shift

# Can You Trust Your Model's Uncertainty? Evaluating Predictive Uncertainty Under Dataset Shift

**Yaniv Ovadia**[*]
Google Research

**Emily Fertig**[*†]
Google Research

**Jie Ren**[†]
Google Research

et al

# nature

# Deep learning and process understanding for data-driven Earth system science

Markus Reichstein ✉, Gustau Camps-Valls, Bjorn Stevens, Martin Jung, Joachim Denzler, Nuno Carvalhais & Prabhat

Conclusions

**(2) Plausibility and interpretability of inferences**
Models should not only be accurate but also credible, incorporating the physics governing the Earth system.

**(3) Uncertainty estimation**
Models should define their confidence and credibility.

# Best practices (I): Hybrid models

BONUS TRACK SLIDE

lightening

machine learning models

physical models

guidance

moreno@dkrz.de

# Best practices (II):
# put your model on diet

Select samples using domain knowledge (lakes, Navier-Stokes examples) ,
use XAI to identify and remove background, spurious correlations (leakage),...



Simonyan et al 2014
https://arxiv.org/pdf/1312.6034v2.pdf

# Best practices (III):
# accuracy is not enough to evaluate the model skills

… the model was very accurate, but in classifying grass (cows example) or it allowed for denser water up (lakes example)...

*"We do not want a correct model, we want understanding"*

Doshi-Velez and Kim, 2017
Towards A Rigorous Science of Interpretable Machine Learning
https://arxiv.org/abs/1702.08608

moreno@dkrz.de

# Best practices (III):
# call a human!

Calculate confidence intervals with uncertainty quantification techniques:

- Conformal Predictors
- MC dropouts
- Deep Ensembles
- Quantile regression
- Gaussian processes (also to propagate distributed input data: https://presentations.copernicus.org/EGU2020/EGU2020-14677_presentation.pdf)

and implement fallbacks if the confidence of the prediction is low.

Maria Navarro: Quantifying uncertainty in Machine Learning predictions | PyData...

PyData • 1.3K views • 6 months ago

conformal predictors
https://github.com/donlnz/nonconformist



Florian Wilhelm: Are you sure about that?! Uncertainty Quantification in AI | PyData...

PyData • 162 views • 1 month ago

Vincent Warmerdam: How to Constrain Artificial Stupidity | PyData London 2019

PyData • 3K views • 6 months ago



GOTO 2018 • Computers are Stupid: Protecting "AI" from Itself • Katharine Jarmul
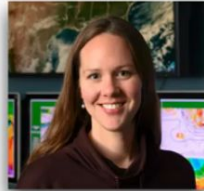
GOTO Conferences ✓ 1.3K views • 12 months ago

# Explainability added value

"Viewing Forced Climate Patterns through an AI Lens", Dec. 11, 2019.