



DSC 5.0 Tutorial: End-to-End ML using h2o in R

Belgrade, November 2019.



About me

- Owner & Consultant @ Logikka
- Senior Data Scientist @ NIS Gazprom Neft
- Board Member @ Data Science Serbia
- Proud ETF alumni

The background features several light blue, semi-transparent abstract shapes. On the left side, there is a large, irregular blob-like shape. In the lower center, there is a solid light blue circle. Another smaller, irregular shape is visible in the upper left quadrant.

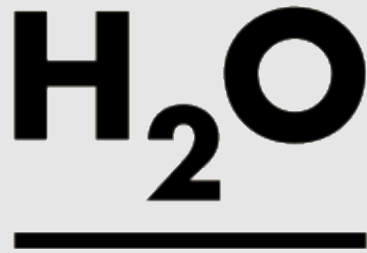
Get tutorial materials:

https://bit.ly/dsc50_h2o_tutorial

WHO
ARE
YOU?



What is h2o?

The logo for H2O, featuring the text "H2O" in a large, bold, black sans-serif font, with a horizontal line underneath the "O".

H₂O

In-memory, distributed
machine learning algorithms

The logo for Sparkling Water, featuring the text "Spark" in a black sans-serif font with an orange star above the "k", followed by a plus sign and "H2O" in a black sans-serif font. Below this is a horizontal line, and then the words "SPARKLING" and "WATER" in a bold, black sans-serif font.

Spark + H₂O
**SPARKLING
WATER**

Integration with Apache
Spark

The logo for H2O4GPU, featuring the text "H2O" in a yellow sans-serif font with a black outline, followed by "4GPU" in a green sans-serif font with a black outline.

H₂O4GPU

Machine Learning on GPUs

100% Open Source



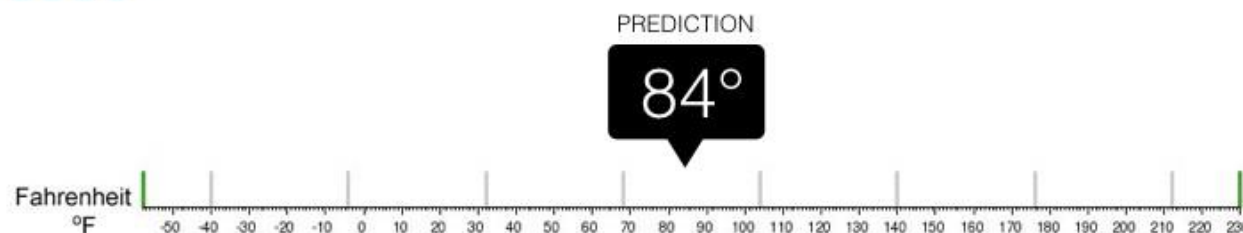
What is Machine Learning?





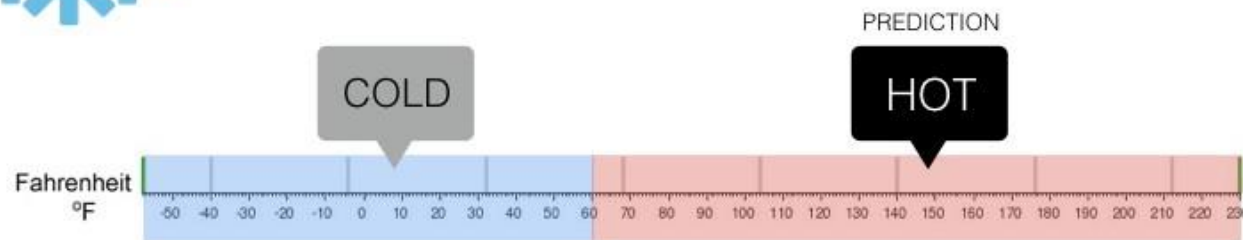
Regression

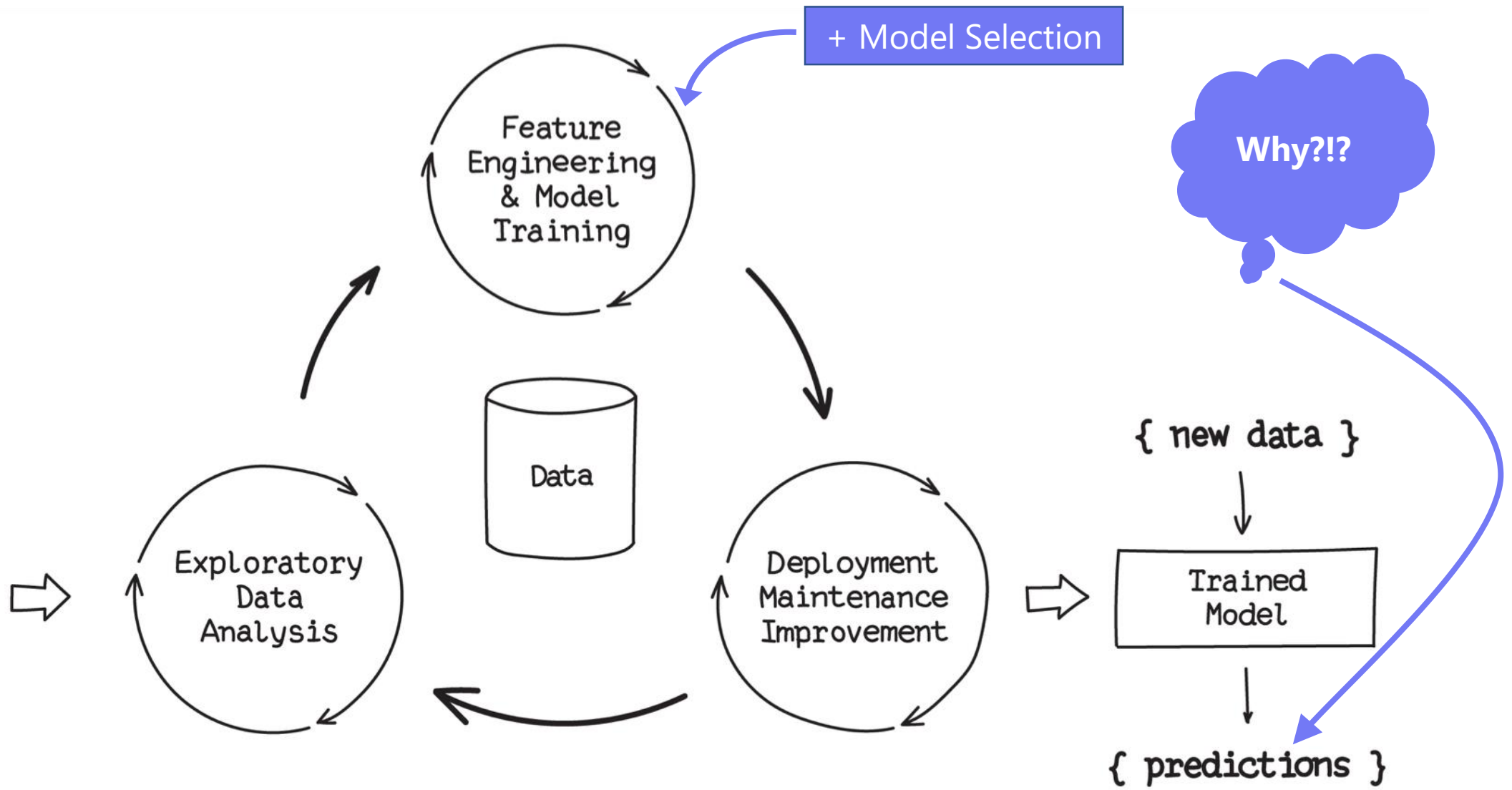
What is the temperature going to be tomorrow?



Classification

Will it be Cold or Hot tomorrow?



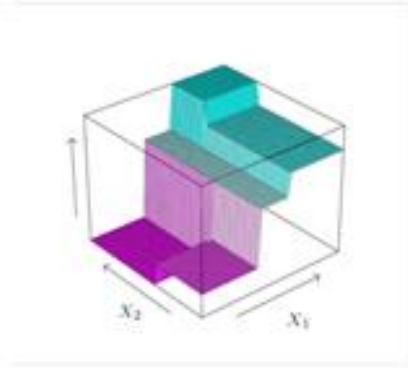


Supported ML Algorithms

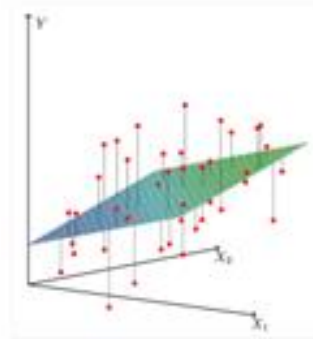
DRF
XRT



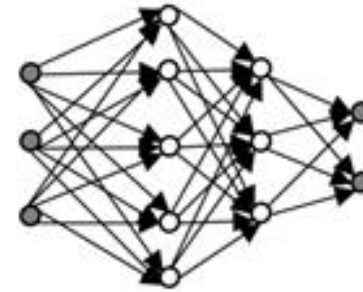
GBM
XGBoost



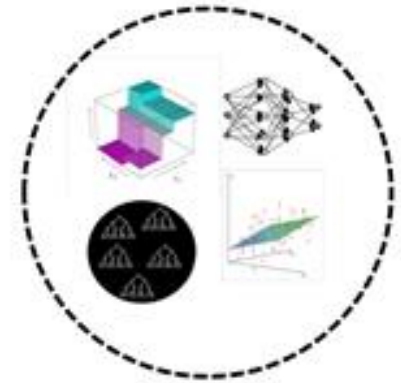
GLM



DNN

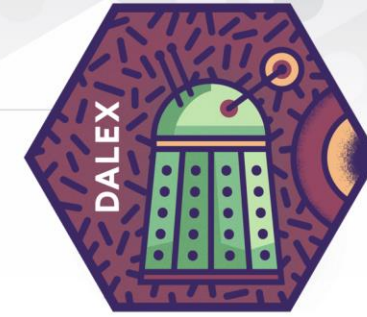


Stacked
Ensemble





DALEX - Descriptive mACHine Learning EXplanations

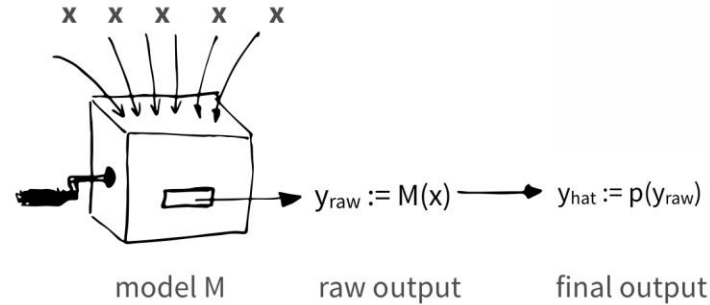


DALEX explains black-box models.
It's a methodology for better
diagnostic of any black-box model.

This approach increases
understanding of a model, increases
trust in model predictions and
allows to further improve the model.
It also allows to compare two or
more models in the scale space

Notation:

- **(x, y)** - pair of input and output data points. x may be anything (data.frame, factors, numbers, text, image), while here we assume that y is numerical or can be transformed to the numerical variable ($x \in X$; $y \in R$).
- **M** - a black box model, $M: X \rightarrow R$. Its output will be denoted as $y_{raw} = M(x)$
- **p** - a link function, transforms raw model output to the same space as y. Useful for classification, while for regression its usually the identity. $p: R \rightarrow R$. Its output will be denoted as $y_{hat} = p(y_{raw})$



`explain(model; data; y; predict_function; trans)`

`prediction_breakdown(explainer, x)`

Prediction explainers shows features that drive model response for a selected observation

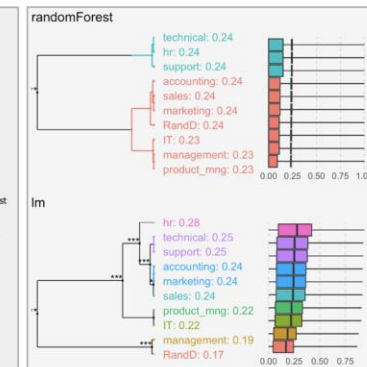
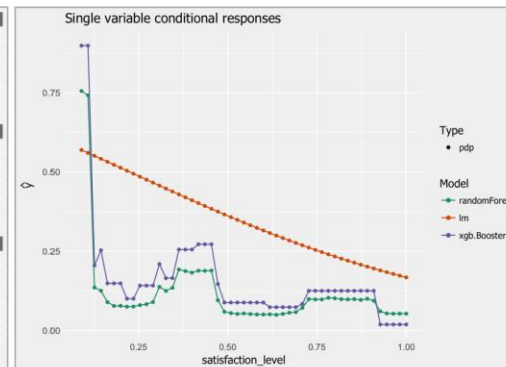
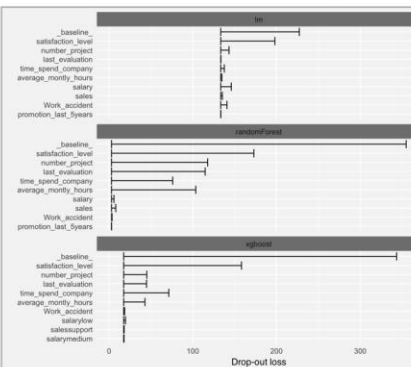
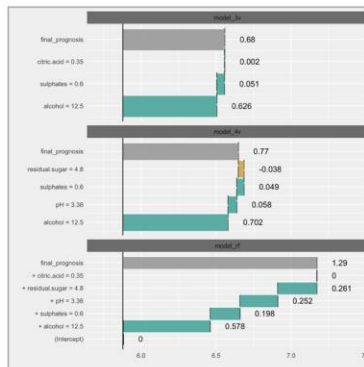
`variable_importance(explainer)`

Variable importance explainers shows the drop in the model loss after permutations of a selected variable.

`variable_response(explainer, variable)`

The `explain()` function creates a wrapper over a black-box model. This wrapper contains all necessary components for further processing.

Single variable explainers show conditional relation between model output and a single variable.



Thank you!

branko@logikka.ai

www.linkedin.com/in/kovacbranko

www.logikka.ai



L O G I K K A