

## Who am I?

## Raphaël Lüthi Machine Learning & Al Lead @ Groupe Mutuel



LinkedIn

2010 EPFL

**Engineering MSE** 

2016

Totaljobs

**Data Scientist** 

2019

**⇔** SBB CFF FFS

**Data Scientist** 

2020

groupemutuel

**Machine Learning & Al Lead** 

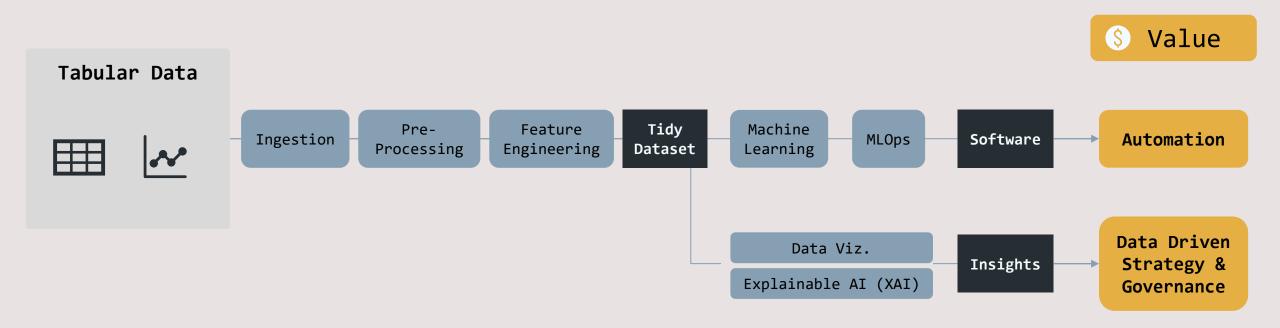
Today



## **Agenda**

- 1. Into: How to Create Value as a Data Scientist?
- 2. The SHAP Explainable AI Method
- 3. Demo: My Workflow on the California Housing Dataset
- 4. Successful Applications in a Business Setting

### **How to Create Value as a Data Scientist?**





## **Agenda**

- 1. Into: How to Create Value as a Data Scientist?
- 2. The SHAP Explainable AI Method
- 3. Demo: My Workflow on the California Housing Dataset
- 4. Successful Applications in a Business Setting

# Data to Insights in Theory

# Using Shap Values to Extract Insights From Data as Proposed by Lundberg et al. 2020

Explainable AI for Trees: From Local Explanations to Global Understanding

Scott M. Lundberg<sup>1</sup>, Gabriel Erion<sup>1,2</sup>, Hugh Chen<sup>1</sup>, Alex DeGrave<sup>1,2</sup>, Jordan M. Prutkin<sup>3</sup>, Bala Nair<sup>4,5</sup>, Ronit Katz<sup>6</sup>, Jonathan Himmelfarb<sup>6</sup>, Nisha Bansal<sup>6</sup>, and Su-In Lee<sup>1,\*</sup>

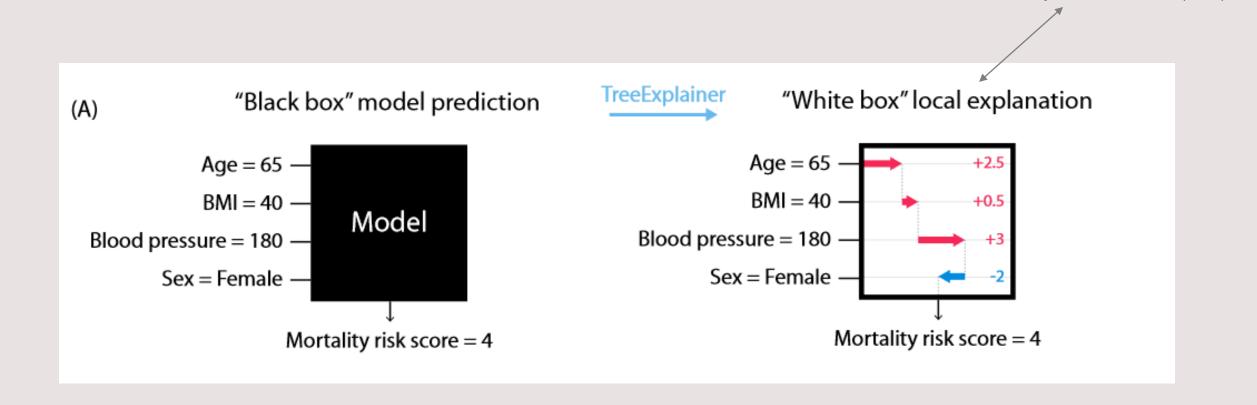
Shapley Additive exPlanations



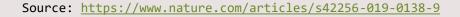


## **Local Explanation with SHAP Values**

Lundberg et al. 2020 :



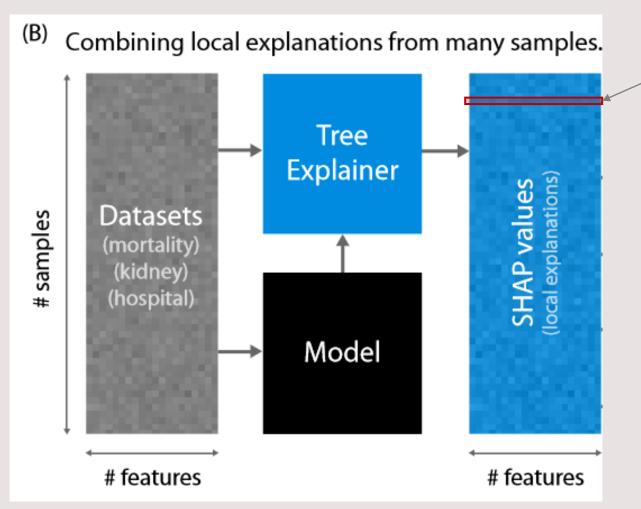
Explainable AI (XAI)

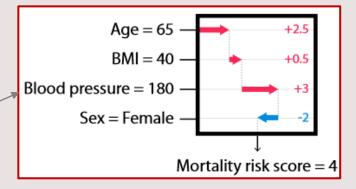




## **Global Explanation with SHAP Values**

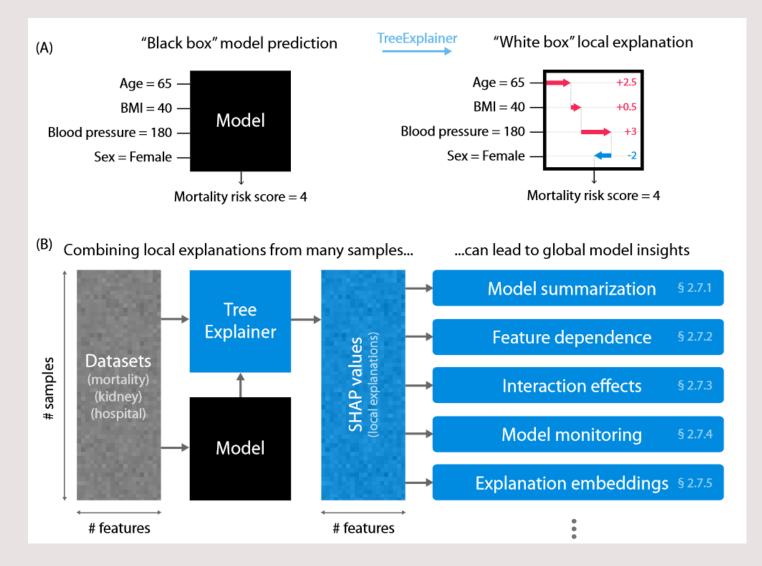
#### Lundberg et al. 2020 :





## **Global Insights with SHAP Values**

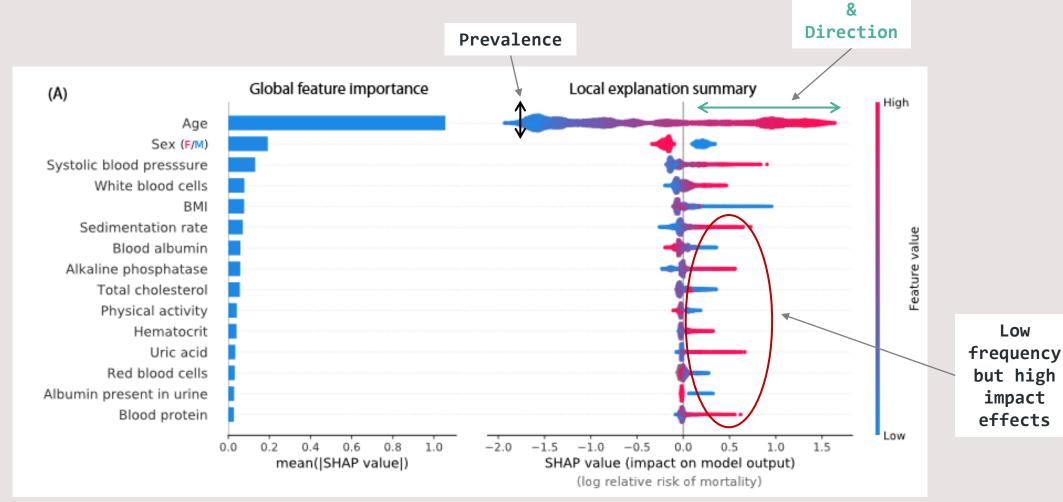
Lundberg et al. 2020 :



Source: <a href="https://www.nature.com/articles/s42256-019-0138-9">https://www.nature.com/articles/s42256-019-0138-9</a>

## **SHAP Summary Plots**

#### Lundberg et al. 2020 :

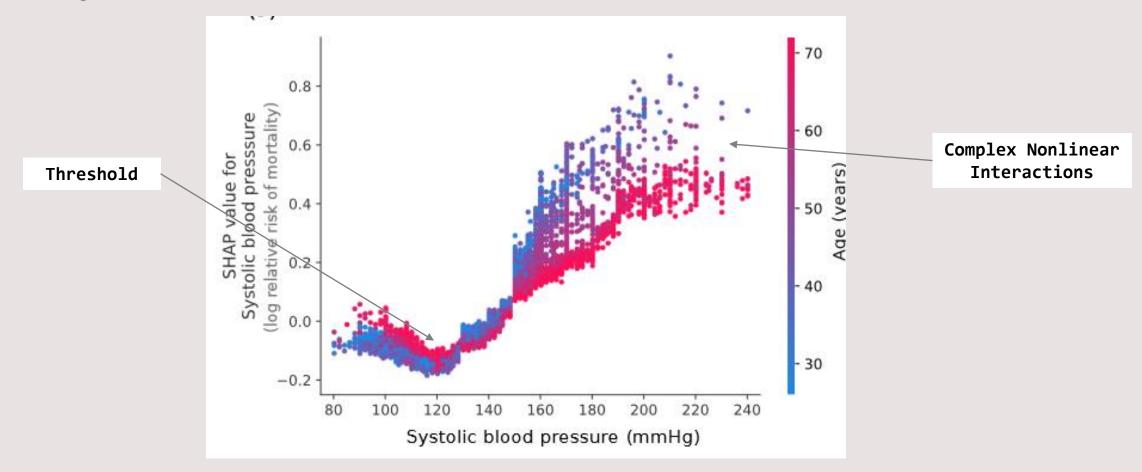


Magnitude



## **SHAP Dependence Plot**

#### Lundberg et al. 2020 :

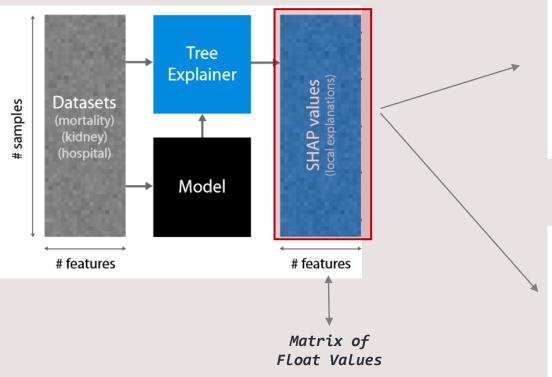




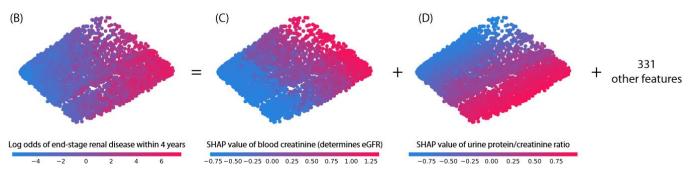
Source: <a href="https://www.nature.com/articles/s42256-019-0138-9">https://www.nature.com/articles/s42256-019-0138-9</a>

## **SHAP Explanation Embedding and Clustering**

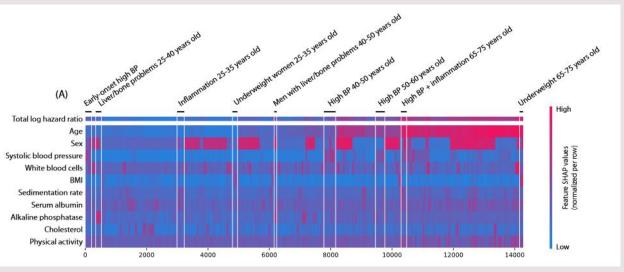
#### Lundberg et al. 2020 :



#### **Explanation Embedding**



#### Explanation Clustering (semi-supervised)





Source: <a href="https://www.nature.com/articles/s42256-019-0138-9">https://www.nature.com/articles/s42256-019-0138-9</a>

## **Agenda**

- 1. Into: How to Create Value as a Data Scientist?
- 2. The SHAP Explainable AI Method
- 3. Demo: My Workflow on the California Housing Dataset
- 4. Successful Applications in a Business Setting

### **Follow Along the DEMO:**

## My Workflow on the California Housing Dataset





**GitHub** 

Google colab

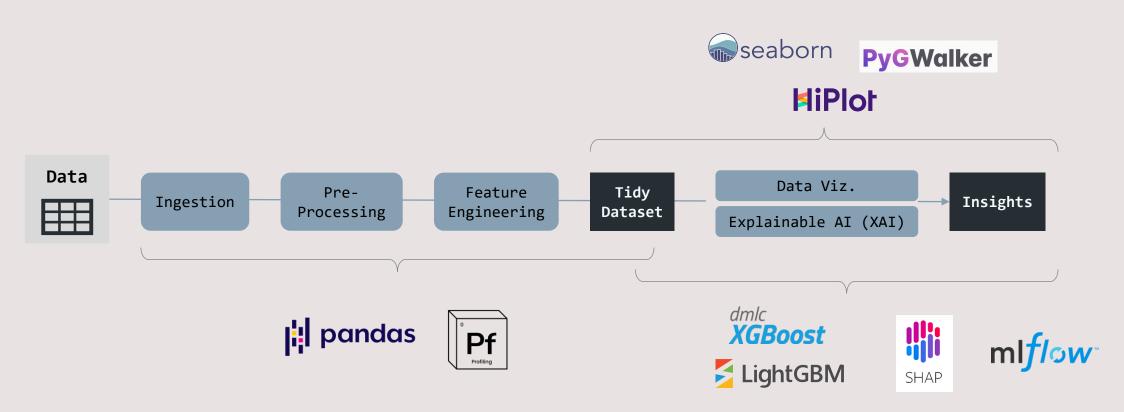
https://github.com/rluthi
/pydata-london-2023

https://colab.research.google.com/github/rluthi
 /pydata-london-2023/blob/main/01 Data-to Insights-Demo PyData-London-23.ipynb



# Data to Insights in Practice

My Battle Tested Workflow Using Open-Source Tools





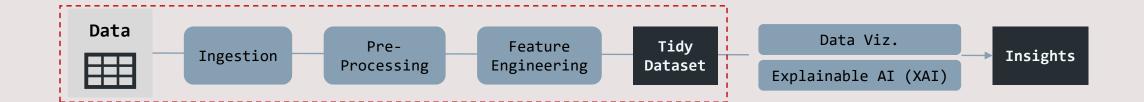
## A Calm Way to Organise Your Data Wrangling



#### pandas pipe

Full tutorial: <a href="mailto:calmcode.io/pandas-pipe/introduction.html">calmcode.io/pandas-pipe/introduction.html</a>

```
data = (
   data raw.pipe(start pipeline)
   .pipe(query, query="MedInc <= 10")
   .pipe(query, query="AveRooms <=10")
   .pipe(query, query=("0.8 <= AveBedrms <=1.4"))
   .pipe(query, query=("Population <= 5000"))
   .pipe(query, query=("AveOccup <= 6"))
   .pipe(query, query=("MedHouseVal <= 5"))</pre>
   .assign(west of lon120=lambda d: d["Longitude"] <= -120)
   .assign(north of lat36=lambda d: d["Latitude"] >= 36)
start pipeline
                                         20,640
                                                  cols:
                                         20,332
                                                                    0.0s | args: query: MedInc <= 10
                                                  cols:
query
                                         20,116
                                                  cols:
                                                           9 | t:
                                                                    0.0s | args: query: AveRooms <=10
query
                                 rows:
                                         19,619
                                                  cols:
                                                                    0.0s | args: query: 0.8 <= AveBedrms <=1.4
                                                           9 | t:
query
                                 rows:
                                         19,317
                                                  cols:
                                                                    0.0s | args: query: Population <= 5000
query
                                 rows:
                                                                    0.0s | args: query: AveOccup <= 6
                                         19,232 | cols:
query
                                         18,562 | cols:
                                                                    0.0s | args: query: MedHouseVal <= 5
query
                                                           9 | t:
```



## **Automate Your EDA**



#### ydata-profiling

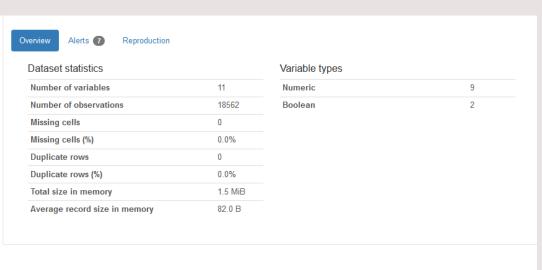
github.com/ydataai/ydata-profiling

from ydata\_profiling import ProfileReport

profile = ProfileReport(df)
profile.to\_file("your\_report.html")

#### **Alternatives**

DataPrep: github.com/sfu-db/dataprep#eda
 Sweetviz: github.com/fbdesignpro/sweetviz



#### Variables

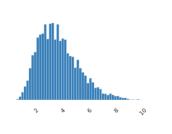


#### MedInc

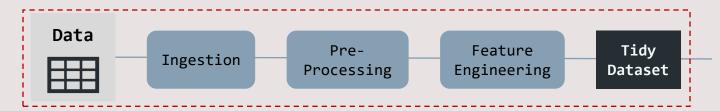
Real number (R)

Distinct	11684
Distinct (%)	62.9%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%
Mean	3.6745204

0.4999
9.9055
0
0.0%
0
0.0%
290.0 KiB



More details



Data Viz.

Explainable AI (XAI)

Insights

## **Explore Your Data Across Many Dimensions**

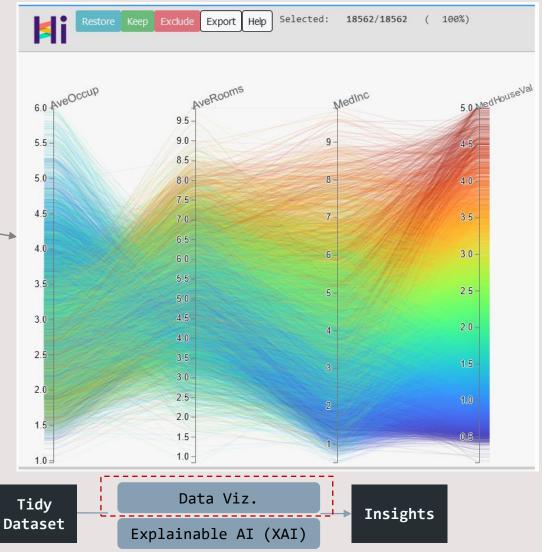
## **HiPlot**

#### HiPlot

facebookresearch.github.io/hiplot

import hiplot

hiplot.Experiment.from dataframe(df).display()





Ingestion

Pre-Processing

Feature Engineering



## **Highlight Patterns With Data Visualisation**

#### PyGWalker

github.com/Kanaries/pygwalker

**PyGWalker** 

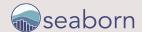
```
import pygwalker

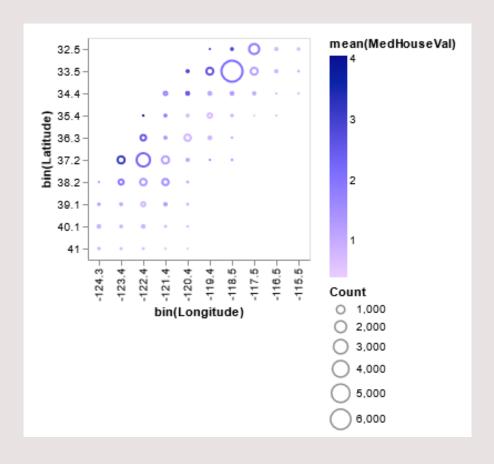
pygwalker.walk(df, axis=1))
```

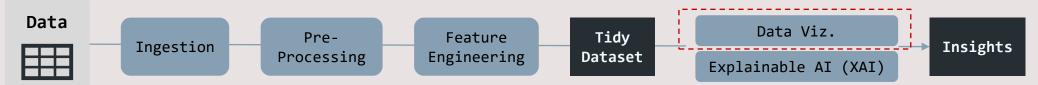
#### seaborn

seaborn.pydata.org/tutorial.html

import seaborn as sns
sns.relplot(data=df, ...)
sns.displot(data=df, ...)
sns.catplot(data=df, ...)

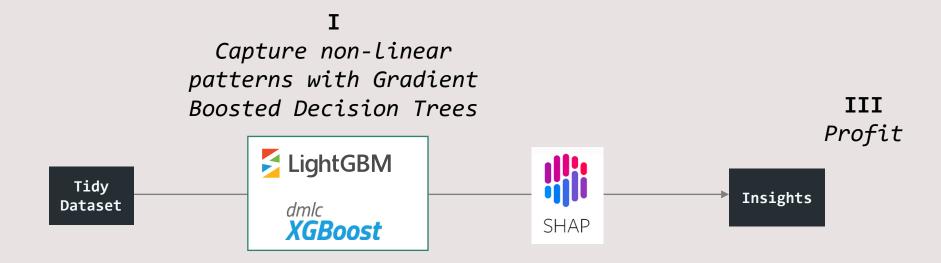






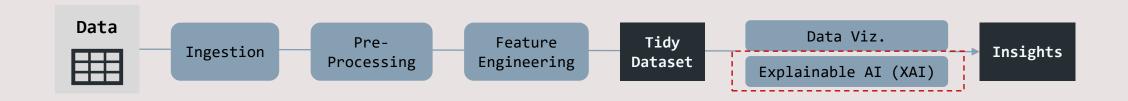


## **Insights from Explainable Al**



II

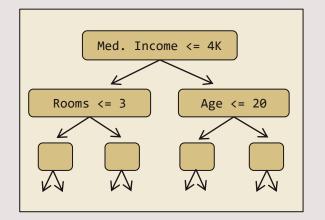
Highlight the patterns
with SHAP



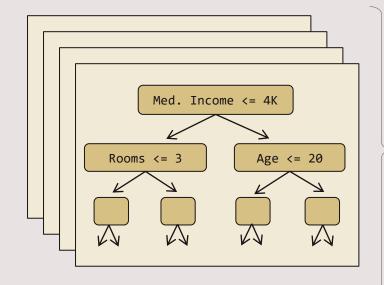
## **Capture Nonlinear Patterns With Gradient**

### **Boosted Decision Trees**

Decision Tree



**Gradient Boosted**Decision Trees

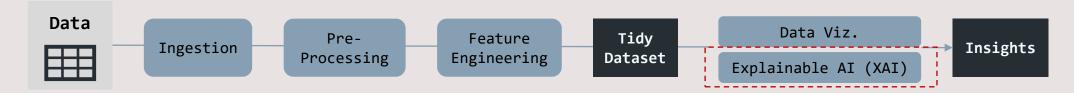






- + Best algorithm for most tabular datasets (benchmarks & Kaggle)
- + Captures nonlinear patterns
- + Great black-box predictors
- Difficult to explain

More info: Pedro Tabacof - Unlocking the Power of Gradient-Boosted Trees (using LightGBM) | PyData London 2022





# Capture Nonlinear Patterns With Gradient Boosted Decision Trees

#### Popular open-source implementations:



#### **LGBM**

LightGBM

lightgbm.readthedocs.io/en/latest



#### XGBoost

eXtreme Gradient Boosting
xgboost.readthedocs.io/en/stable

More info: <a href="Pedro Tabacof">Pedro Tabacof</a> - Unlocking the Power of Gradient-Boosted Trees (using LightGBM) | <a href="PyData London 2022">PyData London 2022</a>

#### Tracking experiments with MLflow



import mlflow

. . .

with mlflow.start run():

mlflow.log\_param("x", 1)

mlflow.log metric("y", 2)

```
import mlflow
mlflow.autolog()
```

- Scikit-learn
- XGBoost
- LightGBM
- Keras
- Gluon
- Statsmodels
- Spark
- Fastai
- Pytorch

```
Data

Ingestion

Pre-
Processing

Pre-
Engineering

Tidy
Data Viz.

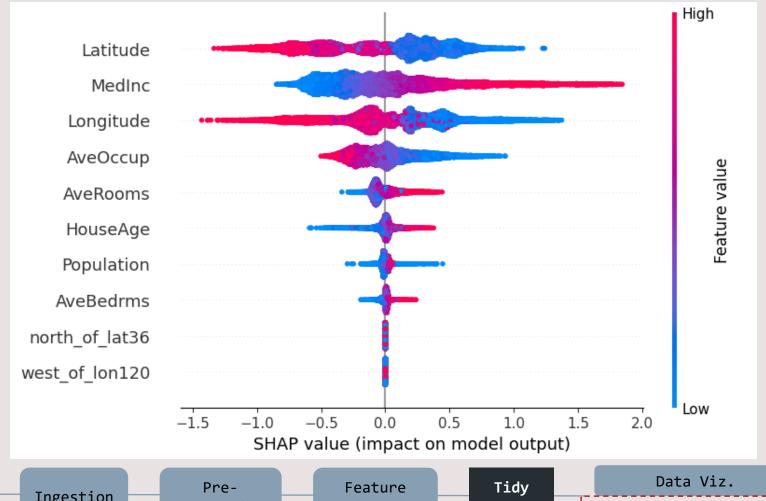
Explainable AI (XAI)

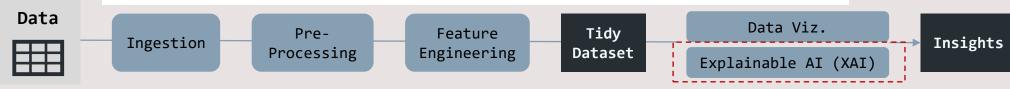
Insights
```



## **Insights from XAI: SHAP Summary Plot**



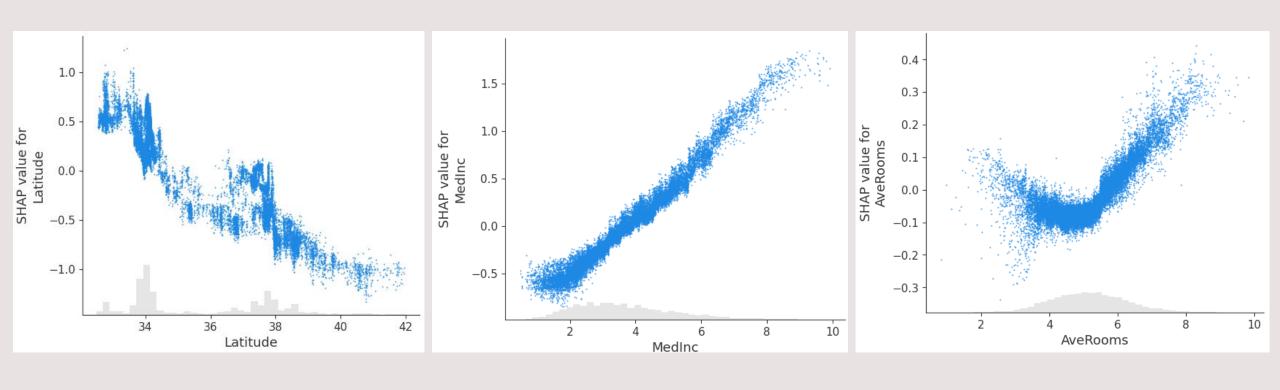


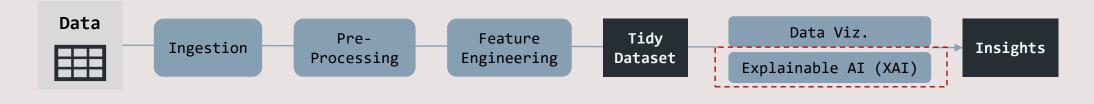




Source: https://shap.readthedocs.io/en/latest/index.html

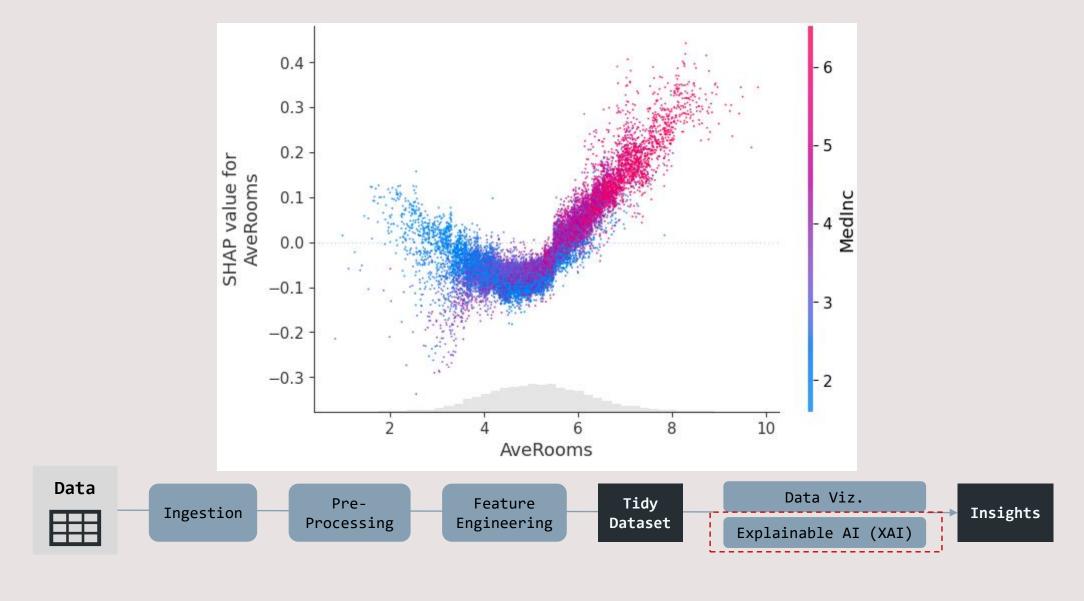
## **Insights from XAI: SHAP Dependence Plots**





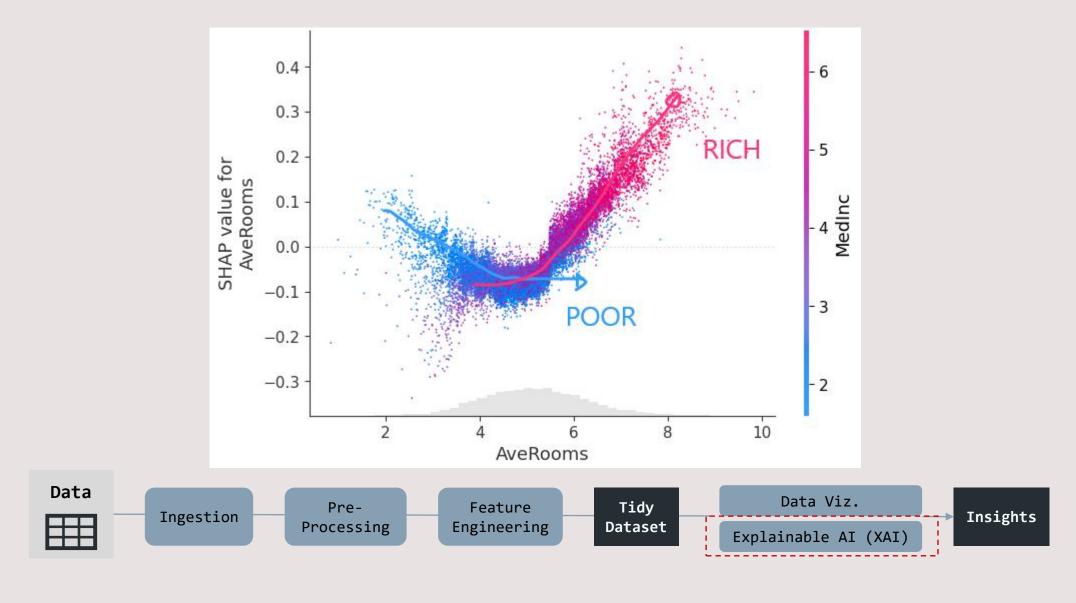


## **Insights from XAI: SHAP Dependence Plots**



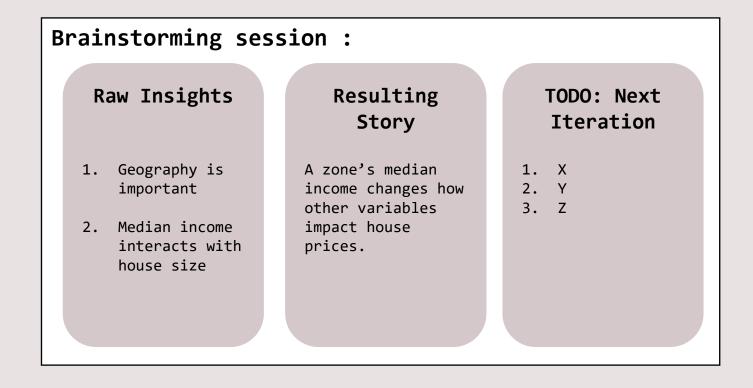


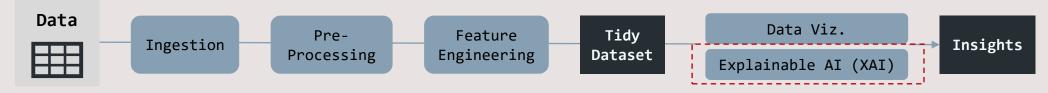
## **Insights from XAI: SHAP Dependence Plots**





# Managing the Complexity & Extracting Business Relevant Insights







## **Agenda**

- 1. Into: How to Create Value as a Data Scientist?
- 2. The SHAP Explainable AI Method
- 3. Demo: My Workflow on the California Housing Dataset
- 4. Successful Applications in a Business Setting

## Successful Applications of XAI as a Product

#### Success stories:

- 1. Which populations were receptive to our marketing campaign ?
- 2. Which populations churn ? And what can we do about them ?
- 3. What are the different user patters on our app?
- 4. Why is this ML model performing suspiciously well?

#### Common business applications:

- Diving into a new dataset and find non-linear patterns
- 2. Segmenting a population by some outcome, we wish to understand or influence
- 3. Reduce the uncertainty of a machine-learning projects



# Data to Insights

Let's Conclude...



## What Should You <u>Take Away</u> From This Talk?

These mature open-source tools will boost your productivity working with tabular data:

- panadas pipe
- Ydata profiling
- HiPlot / PyGWalker / Seaborn
- XGBoost / LightGBM
- MLFlow
- SHAP

You can extract insights by using SHAP (XAI) on Gradient Boosted Decision Trees, resulting in:

- 1. Great model performance
- 2. Transparency in what non-linear patterns are being picking up by the model in the data



## **Thank You for Your Attention!**

## Raphaël Lüthi Machine Learning & Al Lead @ Groupe Mutuel

Get in touch!



LinkedIn

#### Great resources to stay up to date:

- Ian Ozsvald's newsletter:
  - ✓ notanumber.email
- Jesper Dramsch's newsletter:
  - ✓ <u>buttondown.email/jesper</u>
- Vincent Warmerdam's tutorials:
  - ✓ calmcode.io
- Laszlo Sragner's code quality for data science community:
  - ✓ <u>laszlo.substack.com</u>
  - ✓ youtube.com/watch?v=FL6-X1RP7ZE
  - ✓ cq4ds.com
- Avi Chawla's daily blog:
  - √ blog.dailydoseofds.com



## **Thank You for Your Attention!**

