# Causal Machine Learning with DoubleML

Prof. Dr. Martin Spindler, Universität Hamburg & Economic AI

ARIC Lunch Seminar

19.9.2023

# Causal Machine Learning

**ECONOMIC AI**

**Why?** 〉〉 **Causal Machine Learning** 〈〈 **AI**

## Causal Modeling

- Learning causal relationships
- Going beyond correlations
- Pioneers: Pearl, Rubin, Imbens (Nobel Prize 2021)

## Machine Learning

- Learning complex patterns in data
- Correlation based
- Good at forecasting / prediction

# Predictive vs. Causal ML

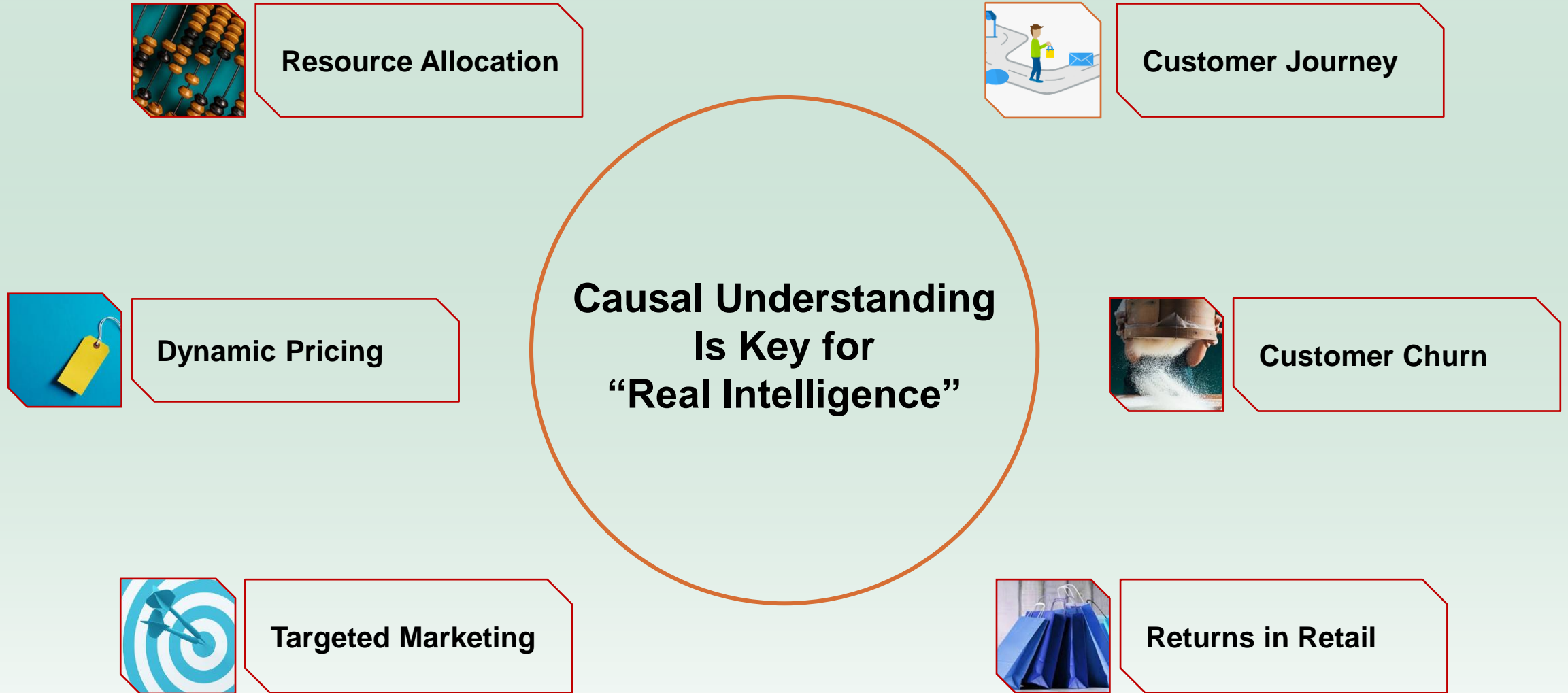| Predictive ML | Causal ML |
|---|---|
| How can we build a good prediction rule, $f(X)$, that uses features $X$ to predict?<br><br>**Example:** Customer Churn<br>    *"How well can we predict whether customers churn?"* | What is the causal effect of a treatment $D$ on an outcome $Y$?<br><br>**Example:** Customer Retention<br>    *"Why do customer churn?"*<br>    *"How can we retain customers?"* |

# Methods

## A/B-Testing / Experimentation

- Getting more popular with tech companies
- Control for covariates to improve precision
- Heterogenous treatment effects

## Structural Economic Models

- Allow for policy evaluation
- Based on economic principles of rational behavior, incentives, etc.
- Allow for competition, strategic effects, etc.

## Hybrid Methods: Instrumental Variables

- Invented in economics but has become popular more broadly
- Used in settings when Randomized Trials are not feasible and/or new policies policy predictions are needed
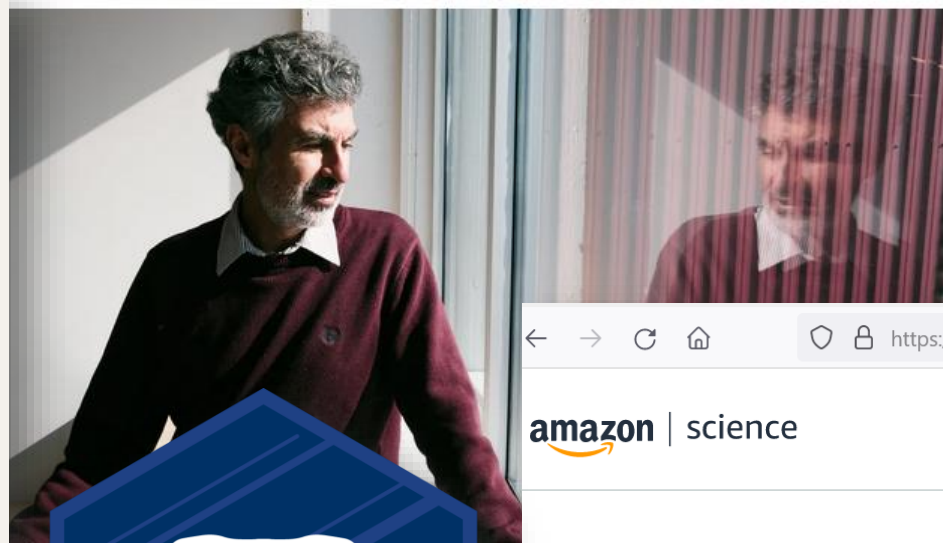- Can handle large dimensions with solid statistical properties

# Causal Machine Learning

doubleml.org

# Application: Randomized Experiments
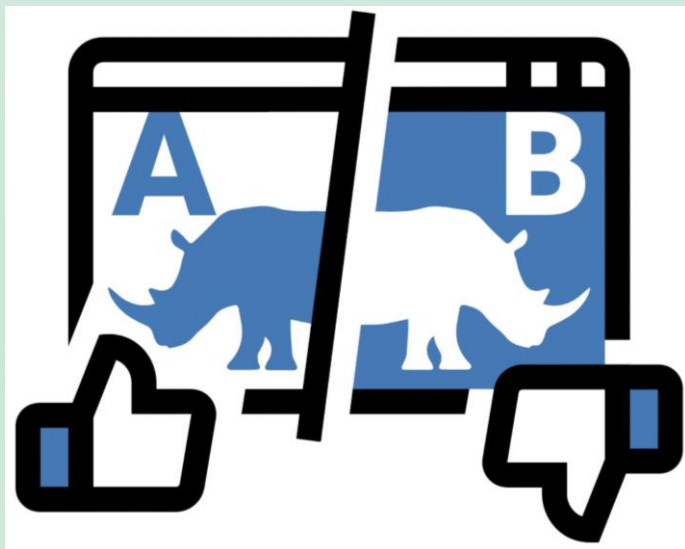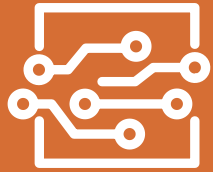
- **General:** What is the effect of a certain variable $D$ on a relevant outcome variable $Y$?

- **Randomized experiments** are a direct way to estimate such effects (assuming they are conducted properly)

## Challenges in practice:

1. No (pure) A/B-testing / experiments possible → observational data

2. A/B test suffers from low power

3. Heterogenous treatment effects

## Solution with DoubleML

1. **Observational study**: Include control variables $X$ which may also impact the variables $Y$ or $D$

2. Include covariates $X$ that help to predict the outcome $Y$ using ML methods

3. Detection of complex treatment effect patterns

# A/B Testing Powered by AI

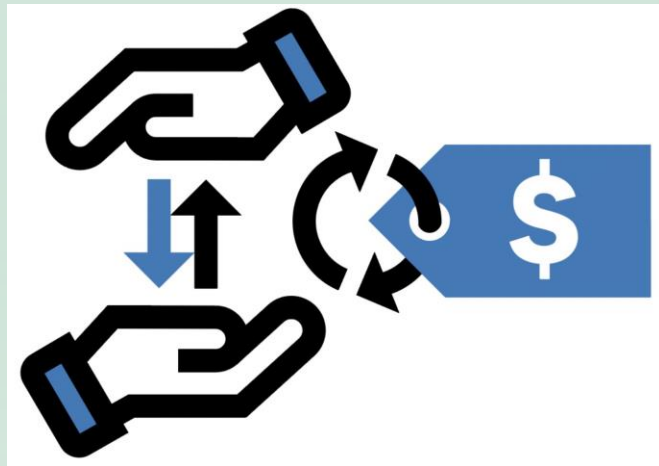**More precise estimation with ML & AI**

**Heterogenous treatment effects & policy optimization**

- „Personalized medicine"
- „Personalized marketing"
- „Dynamic Pricing"

**Adaptive Experiments & Reinforcement Learning**

Price Elasticity of Demand: How does the price impact sales?



- Absolute change in price (EUR 100) and the resulting absolute change in sales (10 million units) can be difficult to interpret

- **Price elasticity of demand**: Percentage change in quantity demanded $D$ when there is a one percent increase in price $P$

$$E_d = \frac{\Delta Q/Q}{\Delta P/P} = \frac{-10/200}{100/1000} = \frac{-0.05}{0.1} = -0.5$$

**Econometric model** for estimating the price elasticity $\theta_0$:

$$\log(Q) = \alpha + \theta_0 \log(P) + X'\beta + \varepsilon,$$

where the vector of controls $X$ can be very **high-dimensional**

## Implemented Extensions

- Different standard causal models (PLM, IRM, IV)
- Simultaneous Inference for Multiple Treatments
- Clustered Standard Errors
- Group Average Treatment Effects (GATEs)
- Conditional Average Treatment Effects (CATEs)
- (Local) Quantile Treatment Effects (QTEs)
- Effects on Conditional Value at Risk (CVaR)

# DoubleML – Models

## Planned Extensions

- DoubleML for difference-in-differences models
- AutoDML
- Sensitivity analysis for omitted variable bias
- Support for unstructured data
- Copula models

# Technical Resources

ECONOMIC AI

## Available Resources

- UAI Tutorial, 2022, available online
- Online Causal Inference Seminar (OCIS), Stanford, 2023
- useR tutorial 2021
- Chamberlain Seminar presentation, 2022, online
- Online documentation (doubleml.org)

# Part II: Use Case
## Development of
## Marketing Promotion Strategy
## For a Software Company

ECONOMIC AI

# Development of Marketing Promotion Strategy

## Initial Situation

- Client: **Company** that **develops and sells enterprise software** to other businesses

- Company used **two types of promotional activities** in the past to increase sales
  - **Free training** for users of the software
  - **Vouchers** to get a discount on purchases

- **Goal**: Company wants to know
  - **Do the promotions actually increase sales?**
  - **Which customer should receive what incentive?**

- **Gold standard**: **Experiment** to test effectiveness
  - ⚡ Company **did not want to do this**

- But: **Data from past transactions** and **promotions available**
  - 💡 **Use Double Machine Learning to estimate effects and derive optimal promotion strategy**

# Confounding Factors Are Present

**Confounding factors are present**

- Customers with **large IT expenses in past received more often incentives** & **correlation between past IT expenses and present sales**

- Potentially other unknown confounders

→ **Need to control for confounders** to get "correct" estimates!

**Goal: Estimate causal effect of promotional activities on sales**

```python
import pandas as pd
import doubleml as dml
from xgboost import XGBRegressor

# Initialize DML data
data_dml = dml.DoubleMLData(
    data,
    y_col='sales',
    d_cols=treatment_vars,
    x_cols=features)

# Instantiate DML model
dml_plr = dml.DoubleMLPLR(
    data_dml,
    ml_l=XGBRegressor(),
    ml_m=XGBRegressor())

# Fit model
dml_plr.fit()
```

$$Y = D'\theta_0 + g_0(X) + \epsilon$$

**Y: Sales**

**D: Promotional activities (& interactions with past revenue)**

- Free training & voucher
- Allow effect of promotions to vary by past revenue (treatment effect heterogeneity!)

**X: Potential confounding variables**

- Large set of customer characteristics, including past IT expenses, # of employees, active worldwide indicator, industry, …

We use a boosted trees algorithm to estimate $g_0(X)$

→ **Fully flexible functional form**

→ **Do not need to worry which variables in X are actually relevant**

# Estimated Effects of Promotional Activities

### Main Effects

€4120**

€3396***

Free Training    Voucher

### Interaction Effects of Promotions and Company's Past Revenue

0,49%**

0,17%**

Free Training    Voucher

**Both types of promotional activities have a positive effect on sales**
- Free training increases sales by €4120 + 0.49% of past revenue
- Voucher increases sales by €3396 + 0.17% of past revenue

Due to confidentiality reasons, the data and estimates are not from the actual project.
However, the insights are qualitatively the same.

Significant at ** 5% level, *** 1% level

# Optimal Marketing Promotion Strategy

**What promotional activity should a customer receive?**

1. Calculate increase in sales for each type of promotional activity, using effect estimates
2. Subtract costs from sales gains
3. **Choose promotion type that yields highest net gain in sales**

**Optimal strategy increases sales by 7.8%** vs. 3.7% under strategy actually observed in data (net of promotion costs)

# Results & Outlook

- **Company recently started to use the developed strategy** to decide how to target customers

- Sales figures since strategy adoption indicate that **strategy yields indeed gains in sales**

- Next steps
  - **Develop additional strategies** for other promotional activities, customer segments, and markets

# More Use Cases for DoubleML



ECONOMIC AI


**Dynamic Pricing**


**A/B Testing**


**Resource Allocation**


**Personalized Marketing**


**Customer Retention**


**… and much more**

Further use cases available upon request!

# References, Resources & Trainings

Prof. Dr. Martin Spindler

Universität Hamburg &
Economic AI

# Online Resources



**DoubleML**

## DoubleML

The Python and R package **DoubleML** provide an implementation of the double / debiased machine learning framework of Chernozhukov et al. (2018). The Python package is built on top of scikit-learn (Pedregosa et al., 2011) and the R package on top of mlr3 and the mlr3 ecosystem (Lang et al., 2019).

### On this page

Main Features
Source code and maintenance
Citation
References

### Getting started

New to **DoubleML**? Then check out how to get started!

### User guide

Want to learn everything about **DoubleML**? Then you should visit our extensive user guide with detailed explanations and further references.

### Workflow

The **DoubleML** workflow demonstrates the typical steps to consider when using **DoubleML** in applied analysis.

ECONOMIC AI

# References – Theory

**Double Machine Learning Approach**

- Chernozhukov, V., Chetverikov, D., Demirer, M., Duo, E., Hansen, C., Newey, W. and Robins, J. (2018), Double/debiased machine learning for treatment and structural parameters. The Econometrics Journal, 21: C1-C68, doi:10.1111/ectj.12097.

- Chernozhukov, V., Hansen, C., Kallus, N., Spindler, M., and Syrgkanis, V. (forthcoming), Applied Causal Inference Powered by ML and AI.

**DoubleML Package for Python and R**

- Bach, P., Chernozhukov, V., Kurz, M. S., and Spindler, M. (2021), DoubleML - An Object-Oriented Implementation of Double Machine Learning in R, arXiv:2103.09603.

- Bach, P., Chernozhukov, V., Kurz, M. S., and Spindler, M. (2022), DoubleML - An Object-Oriented Implementation of Double Machine Learning in Python, Journal of Machine Learning Research, 23(53): 1-6, https://www.jmlr.org/papers/v23/21-0862.html.

ECONOMIC AI

Dynamic Pricing mit Künstlicher Intelligenz

Fallstudie aus dem Ride-Sharing-Markt

Machine learning for financial forecasting, planning and analysis: recent developments and pitfalls

Korrelationen müssen auch kausal sein

# Use Case: Production Optimization

**ECONOMIC AI**



https://proceedings.**mlr**.press/v218/schacht23a.html

## PMLR Proceedings of Machine Learning Research

Volume 218    JMLR    DMLR    TMLR    MLOSS    FAQ    Submission Format    RSS

[edit]

# Causally Learning an Optimal Rework Policy

*Oliver Schacht, Sven Klaassen, Philipp Schwarz, Martin Spindler, Daniel Grunbaum, Sebastian Imhof* Proceedings of The KDD'23 Workshop on Causal Discovery, Prediction and Decision, PMLR 218:3-24, 2023.

## Abstract

# Short Course on Causal Machine Learning

# Outlook – DoubleML Deep

**ECONOMIC AI**

**Economic AI – The Experts for Causal AI Software – Modelling – Trainings**

BOSTON · HAMBURG · HONG KONG · MUNICH

www.economicai.com

**Prof. Dr. Martin Spindler**

spindler@economicai.com

martin.spindler@uni-hamburg.de

linkedin.com/in/martin-spindler

For cooperation and use cases feel free to reach out to us!

**1. Neyman Orthogonality**

The inference is based on a score function $\psi(W; \theta, \eta)$ that satisfies

$$\mathbb{E}[\psi(W; \theta, \eta)] = 0$$

Where $W := (Y, D, X, Z)$ and with $\theta_0$ being the unique solution that obeys the **Neyman orthogonality condition**

$$\partial_\eta \mathbb{E}[\psi(W; \theta, \eta)]\Big|_{\eta=\eta_0} = 0$$

- For many models the Neyman orthogonal score functions are linear in $\theta$

$$\psi(W; \theta, \eta) = \psi_a(W; \eta)\theta + \psi_b(W; \eta)$$

- The estimator $\tilde{\theta}_0$ then takes the form

$$\tilde{\theta}_0 = -(\mathbb{E}_N[\psi_a(W; \eta)])^{-1} \mathbb{E}_N[\psi_b(W; \eta)]$$

PLR example: Orthogonality by including the first-stage regression, i.e., the regression relationship of the treatment variable $D$ and the regressors $X$

Orthogonal score function $\psi(\cdot) = (Y - E[Y|X] - \theta(D - E[D|X]))(D - E[D|X])$

# Neyman Orthogonality

ECONOMIC AI

The two strategies rely on very different moment conditions for identifying and estimating $\boldsymbol{\theta_0}$

$$\mathbb{E}[\psi(W, \theta_0, \eta_0)] = 0$$

**Naive approach**

$$\psi(W, \theta_0, \eta) = (Y - D\theta_0 - g_0(X))D$$

Regression adjustment score

$$\eta = g(X),$$
$$\eta_0 = g_0(X),$$

**FWL partialling out**

$$\psi(W, \theta_0, \eta_0) = ((Y - E[Y|X]) - (D - E[D|X])\theta_0)$$
$$(D - E[D|X])$$

Neyman-orthogonal score (Frisch-Waugh-Lovell)

$$\eta = (g(X), m(X)),$$
$$\eta_0 = (g_0(X), m_0(X)) = (\mathbb{E}[Y \mid X], \mathbb{E}[D \mid X])$$

Both estimators solve the empirical analog of the moment conditions:

$$\frac{1}{n}\sum_{i=1}^{n} \psi(W_i, \theta, \hat{\eta}_0) = 0,$$

where instead of unknown nuisance functions we plug-in their ML-based (hold-out) estimators

# The Key Ingredients of DML

**2. High-Quality Machine Learning Estimators**

The nuisance parameters are estimated with high-quality (fast-enough converging) machine learning methods.

- Different structural assumptions on $\eta_0$ lead to the use of different machine-learning tools for estimating $\eta_0$ (Chernozhukov et al., 2018, Chapter 3)

**3. Sample Splitting**

To avoid the biases arising from overfitting, a form of sample splitting is used at the stage of producing the estimator of the main parameter $\theta_0$.

- Cross-fitting performs well empirically (efficiency gain by switching roles)

# Overview

**DoubleML** provides a general implementation of the Double Machine Learning approach by Chernozhukov et al. (2018) in Python and R

There are also other open-source libraries available for causal machine learning:

- **CausalML** (uber, *https://github.com/uber/causalml*, Chen et al., 2020) - variety of causal ML learners, i.a. with focus on uplift modeling, CATEs and IATEs

- **EconML** (microsoft research, *https://github.com/microsoft/EconML*, Battocchi et al., 2021) – various causal estimators based on machine learning, among others based on double machine learning approach

- ...

**CausalML** and **EconML** have a focus on heterogeneity of treatment effects from their start on

**DoubleML** focuses on implementing the DML approach and its extensions (example: heterogeneity)
- Object-orientated implementation based on orthogonal score
- Extendibility and flexibility

# Building Principles

**Key ingredient and implementation**

- **Orthogonal Score**
  - Object-oriented implementation
  - Exploit common structure being centered around a (linear) score function $\psi(\cdot)$

- **High-quality ML**
  - State-of-the-art ML prediction and tuning methods
  - Provided by `scikit-learn` and `scikit-learn`-like learners

- **Sample Splitting**
  - General implementation of sample splitting

# Why an Object-Orientated Implementation?

Given the components $\psi^a(\cdot)$ & $\psi^b(\cdot)$ of a linear Neyman orthogonal score function $\psi(\cdot)$, a general implementation is possible for

- The estimation of the **orthogonal parameters**

- The computation of the **score** $\psi(W; \theta, \eta)$

- The estimation of **standard errors**

- The computation of **confidence intervals**

- A **multiplier bootstrap** procedure for simultaneous inference

- The **sample splitting** can be implemented in general as well
  - ➢ Implemented in the **abstract base class `DoubleML`**

- The **score components** and the estimation of the **nuisance models** have to be implemented **model-specifically**
  - ➢ Implemented in **model-specific classes** inherited from `DoubleML`

ECONOMIC AI

**DoubleML**

Inference methods based on linear score $\psi$, e.g.,

```
fit()
bootstrap()
confint()
...
```

**partially linear treatment effect**
$$Y = D\theta_0 + g_0(X) + \zeta$$

**binary $D$ & heterogeneous treatment effect**
$$Y = g_0(D, X) + \zeta$$

No —— **Instrumental Variable?** —— Yes

No —— **Instrumental Variable?** —— Yes

**DoubleMLPLR**

Partially linear regression (PLR)

Linear score
$$\psi_a(W;\eta) = -D(D - m(X)),$$
$$\psi_b(W;\eta) = (Y - g(X))(D - m(X)).$$

```
_nuisance_est()
_nuisance_tuning()
...
```

**DoubleMLPLIV**

Partially linear IV regression (PLIV)

Linear score
$$\psi_a(W;\eta) = -(D - r(X))(Z - m(X)),$$
$$\psi_b(W;\eta) = (Y - \ell(X))(Z - m(X)).$$

```
_nuisance_est()
_nuisance_tuning()
...
```

**DoubleMLIRM**

Interactive regression model (IRM)

Linear score
$$\psi_a(W;\eta) = -1,$$
$$\psi_b(W;\eta) = g(1, X) - g(0, X) +$$
$$\frac{D(Y - g(1,X))}{m(X)} - \frac{(1-D)(Y - g(0,X))}{1 - m(x)}.$$

```
_nuisance_est()
_nuisance_tuning()
...
```

**DoubleMLIIVM**

Interactive IV regression model (IIVM)

Linear score
$$\psi_a(W;\eta) = -\frac{D}{p},$$
$$\psi_b(W;\eta) = \frac{D(Y - g(0,X))}{p} - \frac{m(X)(1-D)(Y - g(0,X))}{p(1 - m(X))}.$$

```
_nuisance_est()
_nuisance_tuning()
...
```

# Advantages of the Object-Orientation

**DoubleML** gives the user a **high flexibility** with regard to the specification of DML models:

- Choice of ML methods for approximating the nuisance functions
- Different resampling schemes (repeated cross-fitting)
- DML algorithms DML1 and DML2
- Different Neyman orthogonal score functions

**DoubleML** can be **easily extended:**

- New model classes with appropriate Neyman orthogonal score function can be inherited from DoubleML
- The package features **callables** as score functions which makes it easy to extend existing model classes
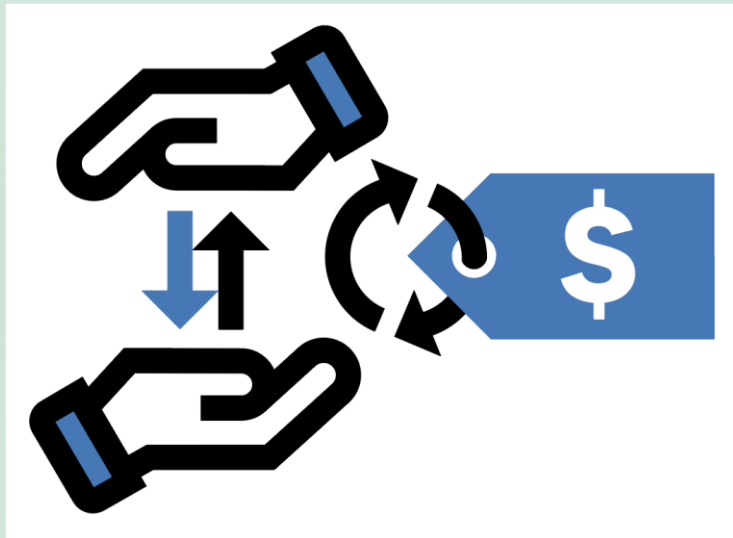- The resampling schemes are customizable in a flexible way

# Installation

Install the latest release via pip or conda, see **installation guide**

```
pip install -U DoubleML
```

```
conda install -c conda-forge doubleml
```

Install development version from GitHub https://github.com/DoubleML/doubleml-for-py

See the Getting Started page of the tutorial website for more information on prerequisites.

# Data Example: Demand Estimation

**Data Source:**
- Data example based on a blogpost by Lars Roemheld (Roemheld, 2021)
- Original real data set publicly available via kaggle, preprocessing notebook available online
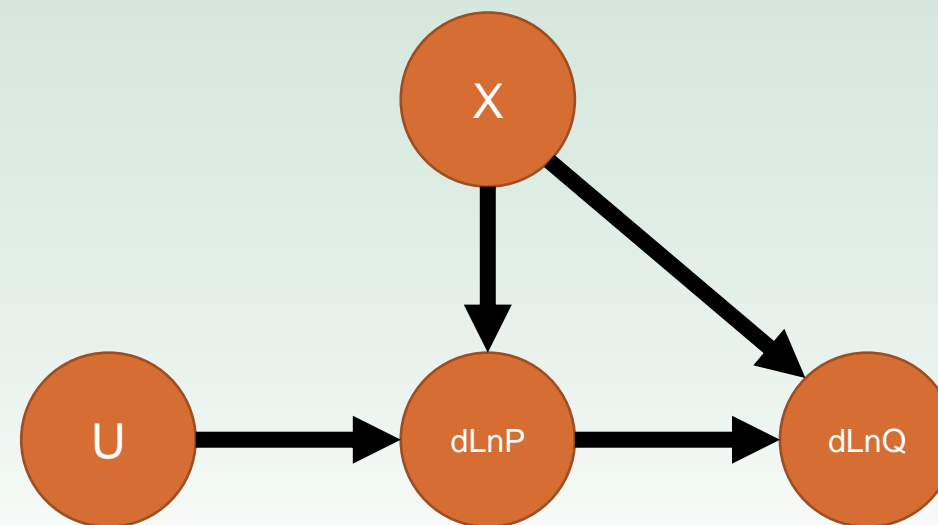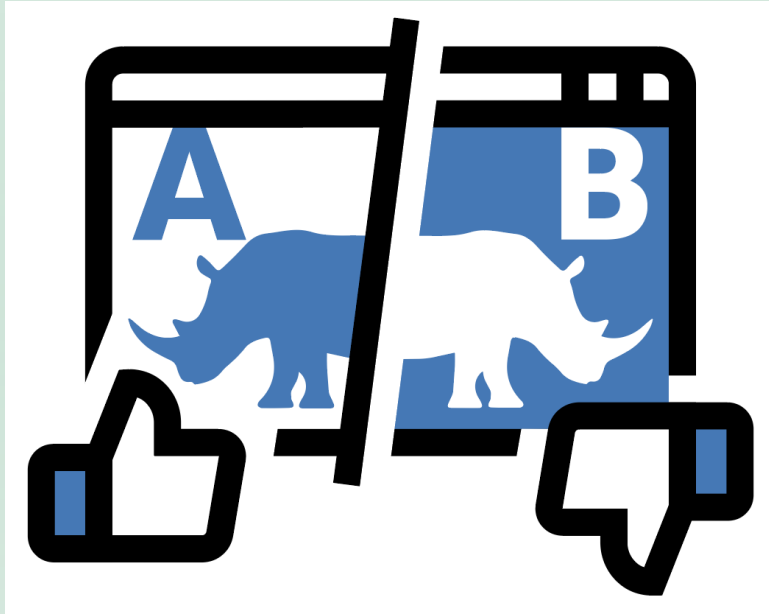
## Causal Problem:

- **Price elasticity of demand**: What is the **effect** of a **price change**, $dLnP$, on **demanded quantity**, $dLnQ$?

- **Observational study**: Flexibly adjust for confounding variables $X$, e.g. product characteristics
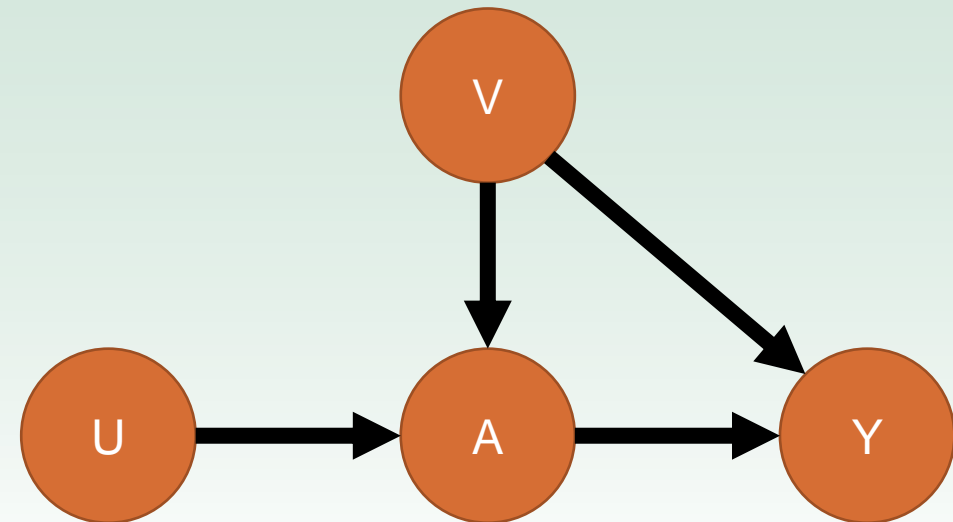
## Causal Diagram (DAG):

**Data Source:**
- Data example based on a randomly chosen DGP created for the 2019 ACIC Data Challenge.

## Causal Problem:

- **Online shop**: What is the effect of a **new ad design** $A$ on sales $Y$ (in $100 )?

- **Observational study**: Necessary to adjust for confounding variables $V$

## Causal Diagram (DAG):

# Online Resources

- The notebook is organized according to the [DoubleML Workflow](#)

- Extensive [User Guide](#) available via [docs.doubleml.org](#)

- [Documentation for the Python API](#) available via [https://docs.doubleml.org/stable/api/api.html](#)

- Paper for the Python package available from [JMLR](#) or [arxiv](#)