



# Active Simulation Data Mining

Mirko Bunse, Amal Saadallah, and Katharina Morik

TU Dortmund, AI Group, 44221 Dortmund, Germany

{firstname.lastname}@tu-dortmund.de

ECML-PKDD 2019  
Interactive Adaptive Learning

### Learning from Simulations

A simulation models how a system  $s$  evolves over time, given the parameters  $\rho$ :

$$\text{Sim}_{\rho}(s_t, \Delta t) = s_{t+\Delta t}$$

Multiple steps make a black-box data generator:

- + Simulated data replaces data from the real system
- + Making predictions is faster than simulating

### Active Sampling

**Goal:** Sample  $\rho$  for maximal learning efficiency

Do so sequentially, monitoring the trained model, and let either  $x \in \rho$  or  $y \in \rho$ , corresponding to the forward/backward distinction.

### Cherenkov Astronomy

Labeled data only through simulation!

**Simulation input:** Label + auxiliary parameters (e.g. energy, direction, ...)

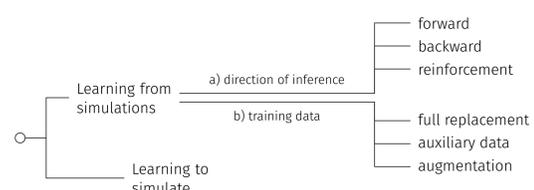
**Simulation output:** Synthetic observations

### Milling

**Simulation input:** Process input parameters

**Simulation output:** Characteristics of the process

Other integrations of machine learning & simulation:



### Forward vs. Backward

One may have to predict the **outcome**  $y$  ...

...or the **cause**  $y'$  of an observation  $x$ .

Example generation is either optimized by active learning (AL) or by active class selection (ACS), depending on the learning scenario.

In particular, the simulation candidates to score are either observations or labels:

$$u_{AL} : \mathcal{X} \rightarrow \mathbb{R}$$

$$u_{ACS} : \mathcal{Y} \rightarrow \mathbb{R}$$

### Sampling of Parameters

The "pure" AL and ACS are *artificially* limited by neglecting the simulation parameters  $\rho$ .

We expect that relevant data can be identified more easily by accounting for all parameters:

$$u : \mathcal{P} \rightarrow \mathbb{R}, \quad \text{where } \begin{cases} \mathcal{X} \subseteq \mathcal{P} & (\text{AL}) \\ \mathcal{Y} \subseteq \mathcal{P} & (\text{ACS}) \end{cases}$$

### Conclusion & Outlook

We distinguish between the *forward* and *backward* learning scenario and propose simulation data mining as a use case for active sampling.

We identify open research issues:

**Transfer Learning:** The simulation may not exactly picture reality. Domain adaptation makes the two domains—simulation and reality—explicit.

**Limits of ACS:** In pure ACS, beating a merely random sampling is hard [5]. By accounting for all parameters  $\rho$ , we hope to find whether this limitation comes from the narrow ACS task itself.

**Data Imbalance in ACS:** Between-class and within-class imbalances may harm ACS strategies.

Want to Collaborate?

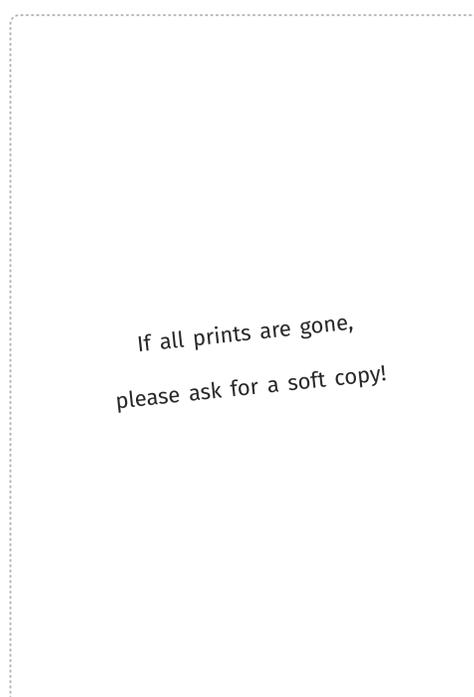


<https://sfb876.tu-dortmund.de/simulation-data-mining>

### References

- [1] C. Bockermann, K. Brügge, J. Buss, A. Egorov, K. Morik, W. Rhode, and T. Ruhe. Online analysis of high-volume data streams in astroparticle physics. In *Proc. of the ECML-PKDD*, pages 100–115. Springer, 2015.
- [2] A. Saadallah, F. Finkeldey, K. Morik, and P. Wiederkehr. Stability prediction in milling processes using a simulation-based machine learning approach. In *51st CIRP Conf. on Manufacturing Systems*. Elsevier, 2018.
- [3] B. Settles. *Active Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2012.
- [4] R. Lomasky, C. E. Brodley, M. Aernecke, D. Walt, and Mark A. Friedl. Active class selection. In *Proc. of the ECML*, pages 640–647. Springer, 2007.
- [5] M. Bunse and K. Morik. What can we expect from active class selection? In *Lernen, Wissen, Daten, Analysen (LWDA) conf. proc.*, 2019. Accepted for publication.

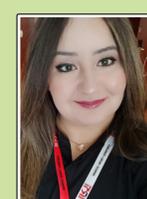
Take this poster with you!



### Mirko Bunse

- Simulation Data Mining
- Deconvolution

...as being applied in Cherenkov astronomy (project SFB 876-C3).



### Amal Saadallah

- Ensemble Methods
- Time Series Analyses

...using sensor data of production PROCESSES (project SFB 876-B3).