# On the Democratization of AI:

## Bridging the Gap with AutoML

---

Candidate:

**Joseph Giovanelli**

Supervisor:
**Matteo Golfarelli**

Coordinator:
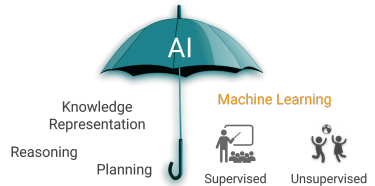**Ilaria Bartolini**

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA
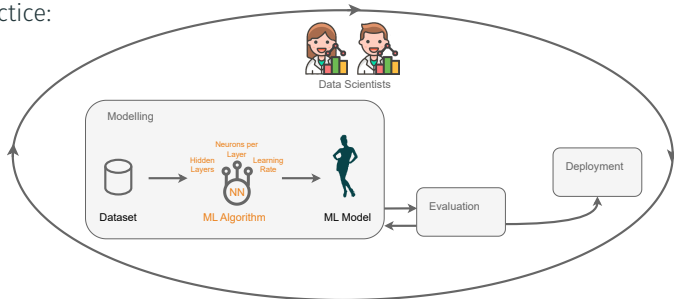
## Machine Learning[1]

A **computer program** is said to learn in some class of **tasks**, with respect to a **performance measure**, if it improves with **experience**.
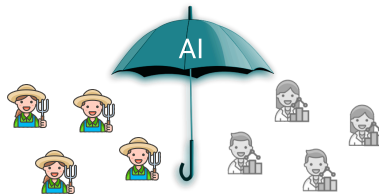


AI

Knowledge Representation

Reasoning

Planning

Machine Learning

Supervised    Unsupervised

In practice:



Data Scientists

Modelling

Dataset

Neurons per Layer

Hidden Layers    Learning Rate

NN

ML Algorithm

ML Model

Evaluation

Deployment

---

[1]Mitchell, T. M. 1997. Machine Learning.
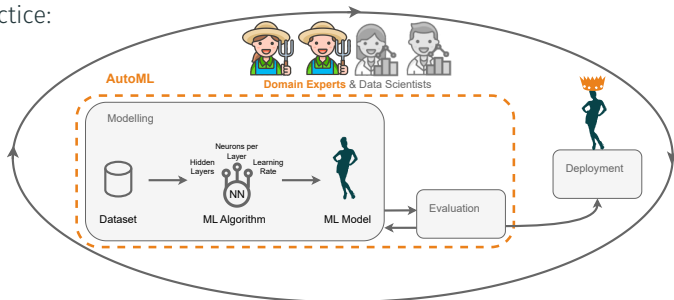
## Democratization of AI[2]

Making AI **accessible to a broader audience**, allowing domain experts to apply it in their own fields.
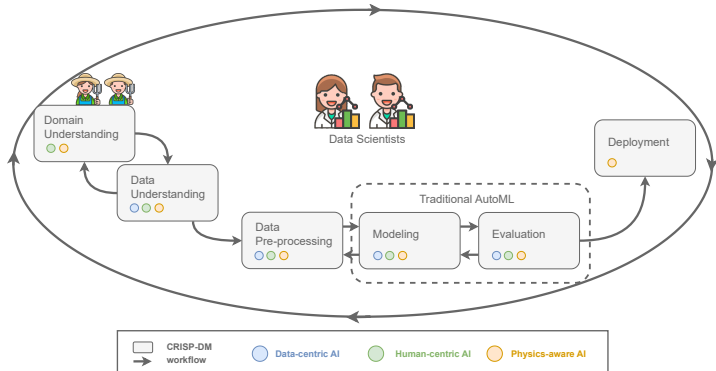


In practice:



---

[2]Thornton, C. et al. 2013. Auto-WEKA: Combined selection and hyperparameter optimization of classification algorithms. In Proceedings of ACM SIGKDD.

Cross-Industry Standard Process for Data Mining (CRISP-DM)
is a process model for dealing with problem complexity.
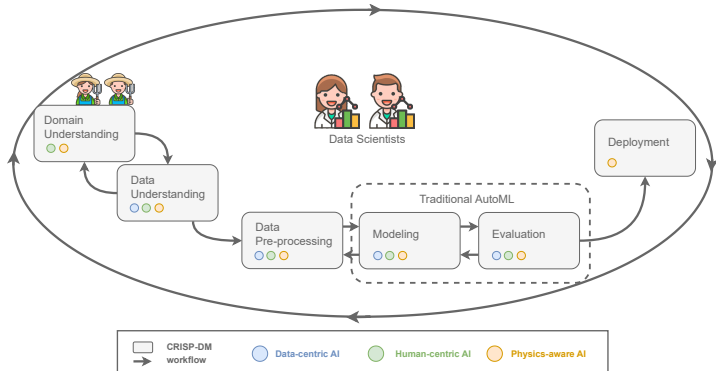


Data-centric AI
systematically engineers data used to build an AI system.

Cross-Industry Standard Process for Data Mining (CRISP-DM)
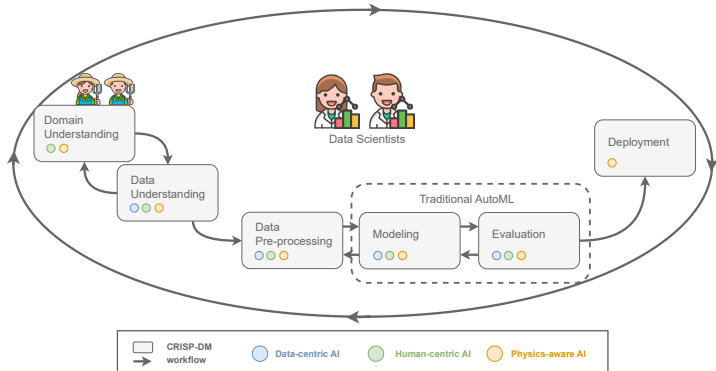is a process model for dealing with problem complexity.



Human-centric AI
aims at complementing, instead of replacing, human intelligence.

## Cross-Industry Standard Process for Data Mining (CRISP-DM)
is a process model for dealing with problem complexity.



## Physics-aware AI
focuses on coupling and enhancing physical simulators with AI.

# Contribution Overview

## Data-centric AI

1. Effective Data Pre-processing Pipelines in Supervised Learning

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

2. Exploring Clustering Pipelines via AutoML and Diversification

$$\mathcal{D} = \{(x_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X}$$

## Human-centric AI

3. Human-centric AutoML via Logic and Argumentat.

$$\boldsymbol{\lambda}^\star \in argmin_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

4. Interactive HPO via Preference Learning

$$\mathcal{L}_1, \ldots, \mathcal{L}_m$$

5. AutoML in the Age of the Large Language Models

## Physics-aware AI

6. Multi-sensor Profiling for Soil-Moisture Monitoring



7. Enhancing Process-Based Models for Soil Moisture Forecasting

**Data-centric AI**  **Human-centric AI**  **Physics-aware AI**

**1.** **Effective Data Pre-processing Pipelines in Supervised Learning [3,4]**

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

**2.** Exploring Clustering Pipelines via AutoML and Diversification

$$\mathcal{D} = \{(x_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X}$$

**3.** Human-centric AutoML via Logic and Argumentat.

$$\boldsymbol{\lambda}^\star \in argmin_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

**4.** Interactive HPO via Preference Learning

$$\mathcal{L}_1, \ldots, \mathcal{L}_m$$

**5.** AutoML in the Age of the Large Language Models

**6.** Multi-sensor Profiling for Soil-Moisture Monitoring



**7.** Enhancing Process-Based Models for Soil Moisture Forecasting



---

[3]<u>Giovanelli J.</u>, Bilalli B., and Abelló A. (2022). Data pre-processing pipeline generation for AutoETL. Information Systems 108 (2022): 101957..

[4]<u>Giovanelli J.</u>, Bilalli B., and Abelló A., et al. (2023). Reproducible experiments for generating pre-processing pipelines for AutoETL. Information Systems (2023): 102314.

# Contribution Overview

**Data-centric AI**

1. Effective Data Pre-processing Pipelines in Supervised Learning

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^N \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

**2. Exploring Clustering Pipelines via AutoML and Diversification** [5]

$$\mathcal{D} = \{(x_i)\}_{i=0}^N \in \mathbb{D} \subset \mathcal{X}$$

**Human-centric AI**

3. Human-centric AutoML via Logic and Argumentat.

$$\boldsymbol{\lambda}^\star \in argmin_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

4. Interactive HPO via Preference Learning

$$\mathcal{L}_1, \ldots, \mathcal{L}_m$$

5. AutoML in the Age of the Large Language Models

**Physics-aware AI**

6. Multi-sensor Profiling for Soil-Moisture Monitoring



7. Enhancing Process-Based Models for Soil Moisture Forecasting



---

[5]Francia M., Giovanelli J., and Golfarelli M. (2024). AutoClues: Exploring Clustering Pipelines via AutoML and Diversification. In Proceedings of Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD). Springer Nature Singapore.

# Contribution Overview

| **Data-centric AI** | **Human-centric AI** | **Physics-aware AI** |

**Data-centric AI**

1. Effective Data
Pre-processing Pipelines
in Supervised Learning

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

2. Exploring Clustering
Pipelines via AutoML and
Diversification

$$\mathcal{D} = \{(x_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X}$$

**Human-centric AI**

3. **Human-centric AutoML
via Logic and Argumentat.** [6]

$$\lambda^\star \in argmin_{\lambda \in \Lambda} \mathcal{L}(A_\lambda(\mathcal{D}_{train}), \mathcal{D}_{val})$$

4. Interactive HPO
via Preference Learning

$$\mathcal{L}_1, \ldots, \mathcal{L}_m$$

5. AutoML in the Age of the
Large Language Models

**Physics-aware AI**

6. Multi-sensor Profiling
for Soil-Moisture Monitoring



7. Enhancing Process-Based
Models for Soil Moisture
Forecasting



---

[6] Francia M., <u>Giovanelli J.</u>, and Pisano G. (2022). **HAMLET: A framework for Human-centered AutoML via Structured Argumentation**. Future Generation Computer Systems 142 (2023): 182-194.

# Contribution Overview

**Data-centric AI**     **Human-centric AI**     **Physics-aware AI**

1. Effective Data Pre-processing Pipelines in Supervised Learning

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

2. Exploring Clustering Pipelines via AutoML and Diversification

$$\mathcal{D} = \{(x_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X}$$

3. Human-centric AutoML via Logic and Argumentat.

$$\boldsymbol{\lambda}^{\star} \in argmin_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

**4. Interactive HPO via Preference Learning [7]**

$$\mathcal{L}_1, \dots, \mathcal{L}_m$$

5. AutoML in the Age of the Large Language Models

6. Multi-sensor Profiling for Soil-Moisture Monitoring

7. Enhancing Process-Based Models for Soil Moisture Forecasting

---

[7] Giovanelli J., Tornede A., Tornede T., and Lindauer M. (2024). Interactive hyperparameter optimization in multi-objective problems via preference learning. In Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 38. No. 11. 2024.

# Contribution Overview

| **Data-centric AI** | **Human-centric AI** | **Physics-aware AI** |
|---|---|---|

1. Effective Data Pre-processing Pipelines in Supervised Learning

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

2. Exploring Clustering Pipelines via AutoML and Diversification

$$\mathcal{D} = \{(x_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X}$$

3. Human-centric AutoML via Logic and Argumentat.

$$\boldsymbol{\lambda}^\star \in argmin_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

4. Interactive HPO via Preference Learning

$$\mathcal{L}_1, \ldots, \mathcal{L}_m$$

5. **AutoML in the Age of the Large Language Models** [8]

6. Multi-sensor Profiling for Soil-Moisture Monitoring

7. Enhancing Process-Based Models for Soil Moisture Forecasting



---

[8]Tornede A., Difan D., Giovanelli J., Mohan A., Ruhkopf T., Segel S., Theodorakopoulos D., Tornede T., Wachsmuth H., and Lindauer M. (2024). **Automl in the age of large language models: Current challenges, future opportunities and risks**. Transaction on Machine Learning Research. ISSN 2835-8856 2024.

# Contribution Overview

**Data-centric AI**

1. Effective Data Pre-processing Pipelines in Supervised Learning

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

2. Exploring Clustering Pipelines via AutoML and Diversification

$$\mathcal{D} = \{(x_i)\}_{i=0}^{N} \in \mathbb{D} \subset \mathcal{X}$$

**Human-centric AI**

3. Human-centric AutoML via Logic and Argumentat.

$$\boldsymbol{\lambda}^\star \in argmin_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

4. Interactive HPO via Preference Learning

$$\mathcal{L}_1, \ldots, \mathcal{L}_m$$

5. AutoML in the Age of the Large Language Models

**Physics-aware AI**

6. **Multi-sensor Profiling for Soil-Moisture Monitoring** [9]

7. Enhancing Process-Based Models for Soil Moisture Forecasting

[9]Francia M., Giovanelli J., and Golfarelli M. (2022). Multi-sensor profiling for precision soil-moisture monitoring. Computers and Electronics in Agriculture. 197 (2022): 106924.

# Contribution Overview

| **Data-centric AI** | **Human-centric AI** | **Physics-aware AI** |
|---|---|---|
| 1. Effective Data Pre-processing Pipelines in Supervised Learning | 3. Human-centric AutoML via Logic and Argumentat. | 6. Multi-sensor Profiling for Soil-Moisture Monitoring |



**Data-centric AI**

1. Effective Data
Pre-processing Pipelines
in Supervised Learning

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^N \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

2. Exploring Clustering
Pipelines via AutoML and
Diversification

$$\mathcal{D} = \{(x_i)\}_{i=0}^N \in \mathbb{D} \subset \mathcal{X}$$

**Human-centric AI**

3. Human-centric AutoML
via Logic and Argumentat.

$$\boldsymbol{\lambda}^\star \in argmin_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

4. Interactive HPO
via Preference Learning

$$\mathcal{L}_1, \dots, \mathcal{L}_m$$

5. AutoML in the Age of the
Large Language Models

**Physics-aware AI**

6. Multi-sensor Profiling
for Soil-Moisture Monitoring

7. **Enhancing Process-Based Models for Soil Moisture Forecasting** [10]

---

[10] Bitelli M., Francia M., Giovanelli J., Golfarelli M., and Tomei F. An Auto-Tuning Three-Dimensional Numerical Model Coupled with Data Assimilation from a Sensor Grid to Forecast Irrigation Demand in Kiwifruit. Submitted to Computers and Electronics in Agriculture.

# Contribution Overview

**Data-centric AI**     **Human-centric AI**     **Physics-aware AI**

1. Effective Data Pre-processing Pipelines in Supervised Learning

$$\mathcal{D} = \{(x_i, y_i)\}_{i=0}^N \in \mathbb{D} \subset \mathcal{X} \times \mathcal{Y}$$

2. Exploring Clustering Pipelines via AutoML and Diversification

$$\mathcal{D} = \{(x_i)\}_{i=0}^N \in \mathbb{D} \subset \mathcal{X}$$

3. Human-centric AutoML via Logic and Argumentat.

$$\lambda^\star \in argmin_{\lambda \in \Lambda} \mathcal{L}(A_\lambda(\mathcal{D}_{train}), \mathcal{D}_{val})$$

**4. Interactive HPO via Preference Learning [7]**

$$\mathcal{L}_1, \dots, \mathcal{L}_m$$

5. AutoML in the Age of the Large Language Models

6. Multi-sensor Profiling for Soil-Moisture Monitoring

7. Enhancing Process-Based Models for Soil Moisture Forecasting

[7] Giovanelli J., Tornede A., Tornede T., and Lindauer M. (2024). Interactive hyperparameter optimization in multi-objective problems via preference learning. In Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 38. No. 11. 2024.

# Interactive Hyperparameter Optimization via Preference Learning in Multi-Objective Problems

# Hyperparameter Optimization (HPO)

**HPO Problem.** Given a machine learning (ML) algorithm $A$ and corresponding hyperparameter space of $\mathbf{\Lambda} = \Lambda_1 \times \cdots \times \Lambda_M$, the goal is to determine the configuration $\mathbf{\lambda}^\star \in \mathbf{\Lambda}$ with optimal loss function $\mathcal{L}$.

$$\mathbf{\lambda}^\star \in \arg\min_{\mathbf{\lambda} \in \mathbf{\Lambda}} \mathcal{L}(A_{\mathbf{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

$\mathcal{L}$ quantifies how well the trained model $h = A_{\mathbf{\lambda}}(\mathcal{D}_{train})$ performs a disjoint split $\mathcal{D}_{val}$.

**Example** NN hyperparameter conf.

$\mathbf{\lambda} \in \mathbf{\Lambda} =$

|  |  |
|---|---|
| *learning rate:* | 0.05 |
| *hidden layers:* | 5 |
| *neurons per layer:* | 256 |



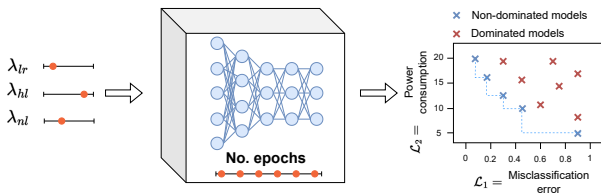**✕** Model  $h = A_{\mathbf{\lambda}}(\mathcal{D}_{train})$

**SOTA.** Bayesian Optimization (BO)[11] drives the exploration toward new **promising configurations** via a surrogate trained on past evaluations.

---

[11]Brochu E., Vlad M. Cora, et al. **A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning.** (2010).

# Hyperparameter Optimization (HPO)

**HPO Problem.** Given a machine learning (ML) algorithm $A$ and corresponding hyperparameter space of $\boldsymbol{\Lambda} = \Lambda_1 \times \cdots \times \Lambda_M$, the goal is to determine the configuration $\boldsymbol{\lambda}^\star \in \boldsymbol{\Lambda}$ with optimal loss function $\mathcal{L}$.

$$\boldsymbol{\lambda}^\star \in \arg\min_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \mathcal{L}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}), \mathcal{D}_{val})$$

$\mathcal{L}$ quantifies how well the trained model $h = A_{\boldsymbol{\lambda}}(\mathcal{D}_{train})$ performs a disjoint split $\mathcal{D}_{val}$.

**Example** Best NN hyperparameter conf.

$\boldsymbol{\lambda}^\star \in \boldsymbol{\Lambda} =$
| | |
|---|---|
| `learning rate:` | 0.01 |
| `hidden layers:` | 10 |
| `neurons per layer:` | 256 |



**SOTA.** Bayesian Optimization (BO)[11] drives the exploration toward new **promising configurations** via a surrogate trained on past evaluations.

[11]Brochu E., Vlad M. Cora, et al. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. (2010).

# Multi-Objective Machine Learning (MO-ML)

MO-ML algorithms. When optimizing multiple objectives $\mathcal{L}_1, \ldots, \mathcal{L}_M$, MO-ML algorithms $A_{\boldsymbol{\lambda}}(\mathcal{D}_{train})$ return a **Pareto front** $P_{\mathcal{D}_{val}}(\mathcal{H})$.



Quality Indicators. quantify the goodness of the Pareto front by measuring specific characteristics —e.g., hypervolume (HV)[13], maximum spread (MS)[13], spacing (SP)[14], <u>closeness to reference point (R2)[15]</u>.

[13]Zitler E., Thiele L. **Multiobjective evolutionary algorithms: a comparative case study and the strength Pareto approach.** IEEE transactions on Evolutionary Computation, 1999.
[14]Schott J. R. **Fault tolerant design using single and multicriteria genetic algorithm optimization.** 1995. PhD Thesis. Massachusetts Institute of Technology.
[15]Hansen M. P., Jaszkiewicz A. **Evaluating the quality of approximations to the non-dominated set**. Copenhagen, Denmark: IMM, Technical University of Denmark, 1994.

6

**Challenge.** Choosing the **quality indicator** leading to a Pareto front which has a **desired shape** requires deep expert knowledge.

**Approach.** Learning a quality indicator via user preferences

1. **Preliminary Sampling**
   Collect Paretos $\mathcal{P} = \{P_i | P_i = P_{\mathcal{D}_{val}}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train})) : \boldsymbol{\lambda} \in \boldsymbol{\Lambda}\}$

2. **Interactive Preference Learning**
   a. Build preference dataset $\mathcal{U} = \{P_{i,1} \succ P_{i,2}\}_{i=1}^{U}$
   b. Learn utility function $u : \mathcal{P} \to \mathbb{R}$

3. **Utility-driven HPO**

   Solve HPO problem $\boldsymbol{\lambda}^{\star} \in \arg\min u(P_{\boldsymbol{\lambda}}) : P_{\boldsymbol{\lambda}} = A_{\boldsymbol{\lambda}}(\mathcal{D}_{val})$

# 1. Preliminary Sampling



Example $\boldsymbol{\lambda_1} \in \boldsymbol{\Lambda}$

Example $P_{\mathcal{D}_{val}}(A_{\boldsymbol{\lambda}}(\mathcal{D}_{train}))$

NN hyperparameter configuration:

| | |
|---|---|
| *learning rate:* | 0.01 |
| *hidden layers:* | 10 |
| *neurons per layer:* | 256 |

# 1. Preliminary Sampling



Example $\lambda_2 \in \Lambda$

Example $P_{\mathcal{D}_{val}}(A_{\lambda}(\mathcal{D}_{train}))$

NN hyperparameter configuration:

| | |
|---|---|
| *learning rate:* | 0.05 |
| *hidden layers:* | 5 |
| *neurons per layer:* | 256 |

Example $\mathcal{U} = \{P_{i,1} \succ P_{i,2}\}_{i=1}^{U}$

$$P_{i,1} \qquad\qquad \succ \qquad\qquad P_{i,2}$$
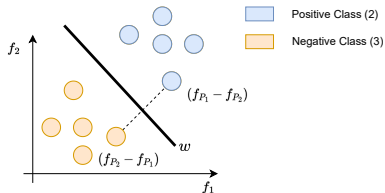
Example
$f_{P_{\mathcal{D}_{val}(H)}} \in \mathbb{R}^d$

**Example** Building $u : \mathcal{P} \to \mathbb{R}$ through RankSVM[16]

$$P_1 \succ P_2 \Leftrightarrow \vec{w}^\mathsf{T} \vec{f}_{P_1} > \vec{w}^\mathsf{T} \vec{f}_{P_2} \quad (1)$$
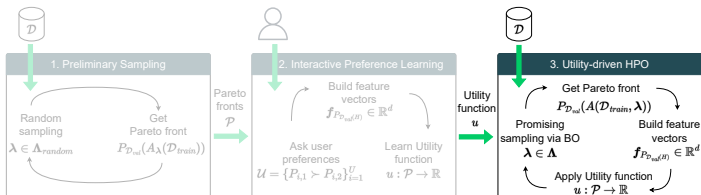
$$\Leftrightarrow \vec{w}^\mathsf{T} \left( \vec{f}_{P_1} - \vec{f}_{P_2} \right) > 0 \quad (2)$$

$$\Leftrightarrow \vec{w}^\mathsf{T} \left( \vec{f}_{P_2} - \vec{f}_{P_1} \right) < 0 . \quad (3)$$



---

[16] Joachims T. **Optimizing search engines using clickthrough data.** Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining. (2002).

**Example.**
Promising sampling
$\boldsymbol{\lambda} \in \boldsymbol{\Lambda} \qquad \longrightarrow$

$learning\ rate:$ 0.001
$hidden\ layers:$ 5
$neurons\ layer:$ 256

**Example.**
Get Pareto and build feature vector
$P_{\mathcal{D}_{val}}(A(\mathcal{D}_{train}, \boldsymbol{\lambda})) \qquad \longrightarrow \qquad \boldsymbol{f}_{P_{\mathcal{D}_{val}}}$



**Example.**
Apply utility
$u(\boldsymbol{f}_{P_{\mathcal{D}_{val}}}) \qquad \longrightarrow$

$w^T \cdot \boldsymbol{f}_{P_{\mathcal{D}_{val}}(H, \boldsymbol{\lambda})}$

**Preference-Based** (PB): HPO process driven by the utility function trained with the indicator in the row

**Indicator-Based** (IB): HPO process driven by the indicator in the column

LCBench[17]:

- funnel-shaped **MLP** from **Auto-pytorch**;
- 35 datasets from **OpenML CC-18** suite.

| PB\IB | HV | SP | MS | R2 |
|-------|--------|---------|---------|---------|
| HV | 98.70% | 146.15% | 146.15% | 98.70% |
| SP | 300.00% | 100.00% | 400.00% | 400.00% |
| MS | 321.05% | 321.05% | 93.85% | 265.22% |
| R2 | 95.65% | 204.35% | 195.65% | 95.65% |

$\longrightarrow$ PB performs **better or equal** in 11/16 cases;
$\longrightarrow$ IB performs **slightly better** in only 5/16 cases.

---

[17] Zimmer L., Lindauer M, and Hutter F. **Auto-pytorch: Multi-fidelity metalearning for efficient and robust autodl.** IEEE transactions on pattern analysis and machine intelligence 43.9 (2021).
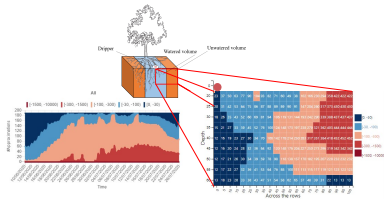
# Conclusions and Future Works

## Main Contributions
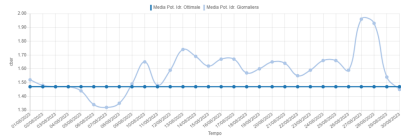
- Develop effective data pipelines for **supervised** and **unsupervised learning**.

- Propose a HAMLET: a human-centric AutoML framework allowing **explainability** and **interactive interventions**;

## Future Works

- Provide insight on data bias in pre-processing for **fairness** through HAMLET.

- Integrate HAMLET with **multi-objective** and **cross-cutting constraints** (e.g., ethical, legal) to allow fairness interventions.

# Physics-aware AI





## Main Contributions

- Integration of physical models with AI through AutoML for **monitoring** and **forecasting tasks**.

- During the whole campaign:
  Water saving: 44%
  Vine productivity: unaffected
  Fruit quality: increased (+1 brix)

## Future works:

- Integration of a smart-irrigation algorithm based on **control theory**, dynamically adjusting the water plan.

- Application of **transfer learning** with a pre-trained model to transfer the knowledge from different conditions, supporting a wider range of crops.

Thanks for the attention :)