

AutoML

A state-of-the-art overview

Joseph Giovanelli

j.giovanelli@unibo.it

Alma Mater Studiorum • University of Bologna

BIG • Business Intelligence Group



Ph.D. Candidate

in Computer Science and Engineering

Main research field: **AutoML**



UPC - Barcelona

Meta-learning, Data

Pre-processing

LUH|AI - AutoML Hannover

Multi-objective, Preference

Learning



Research & Development

projects on BI, Big Data, Data Mining

Table of contents

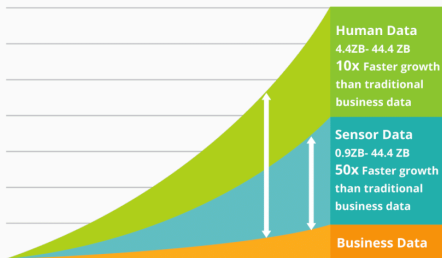
1. Introduction
2. Building blocks
3. State of the art
4. Human-centered AutoML

Introduction

The data growth

It has been reported that 2.5 quintillion bytes of data is being created everyday

The 90% of stored data in the world, has been generated in the past two years only ¹



¹Forbes: How Much Data Do We Create Every Day? May 21, 2018

The sexiest job of the 21st century

The **Data scientist** has become one of the most sought figure

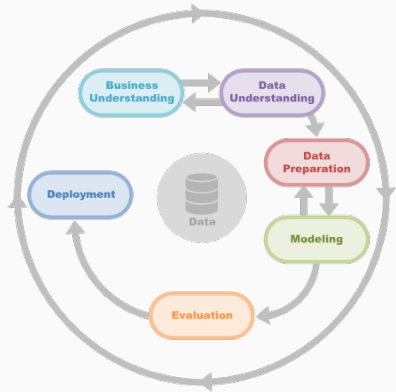
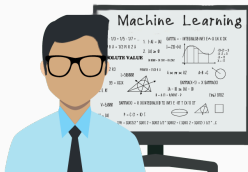
proschool
An IBM Initiative

DATA SCIENTIST MUST-HAVE SKILLS

- MATH & STATISTICS**
 - Machine Learning
 - Statistical Modelling
 - Exploratory Analysis
 - Clustering
 - Regression Analysis
- PROGRAMMING & DATABASE**
 - Computer Science Fundamentals
 - Database Management System
 - Data Visualization
 - Python
 - Big Data
- DOMAIN KNOWLEDGE & SOFT SKILLS**
 - Inclination towards business operations
 - Keen on working with data
 - Problem solver
 - Strategic, proactive, and cooperative
 - Interested in hacking
- COMMUNICATION & VISUALIZATION**
 - Storytelling skills
 - Convert data-based insights into decisions
 - Collaborative with Sr. Management
 - Knowledge of tools like Tableau
 - Visual art design

The role of Machine Learning in Data Science

Data scientists use the **Machine Learning** toolbox to solve real-cases problems



The need

Data Scientists do not scale: ²

- the **increasingly growing size of data** overcomes their availability
- the **numerous skills expected** (IT, mathematics, statistics, business, cooperation) make it difficult to increase their number



More and more **non-experts use data mining tools**

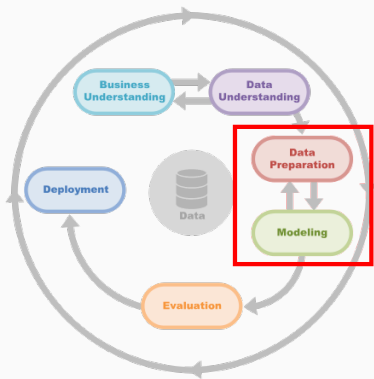


Off the shelf solutions are required to assist them

²Harvard Business Review: Data Scientists Don't Scale, May 22, 2015

AutoML definition

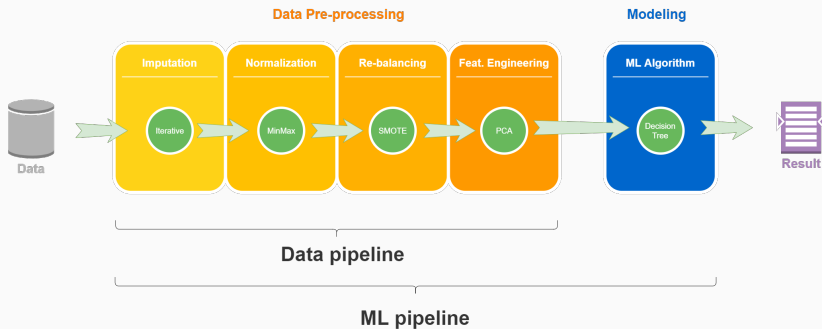
Automated Machine Learning is the process of automating the process of applying **Machine Learning**



Data scientists can spend less tedious time on finding parameters/hyperparameters, and focus on the analysis

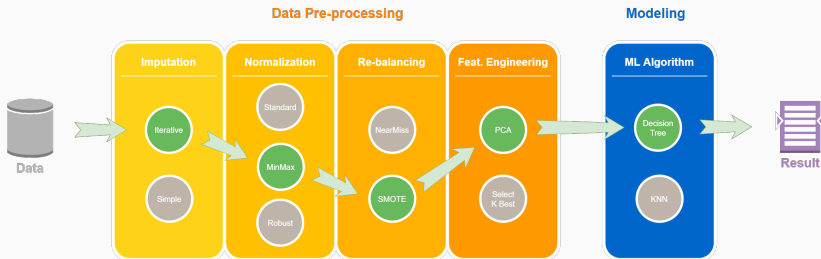
AutoML outcome

AutoML aims to find a **ML pipeline**



AutoML outcome

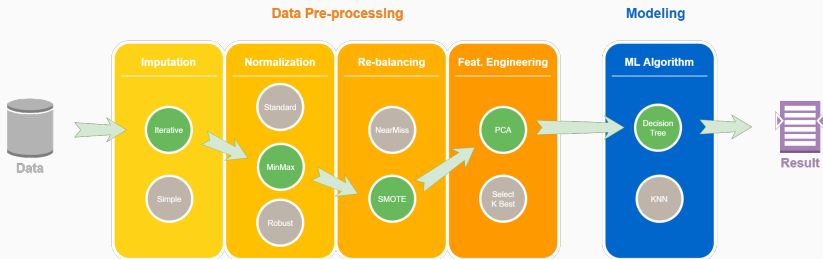
AutoML smartly explores huge search spaces.



- A **data pipeline** consists of a **sequence of transformations**
- Each **transformation** can be instantiated from a pool of **operators**
- Each **operator** has several **parameters**
- Each **parameter** has its own **search space**

AutoML outcome

AutoML smartly explores huge search spaces.

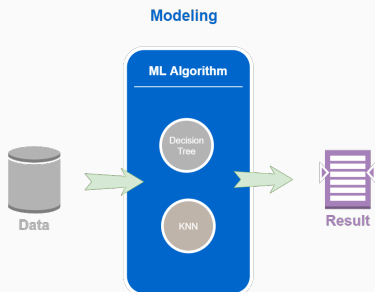


- The **modeling phase** involves the instantiation of a **algorithm** from a specific
- Each **algorithm** has several **hyper-parameters**
- Each **hyper-parameter** has its own **search space**

Building blocks

Auto-WEKA: the CASH problem

Auto-Weka introduces the Combined Algorithm Selection and Hyper-parameter optimization problem (CASH)³



DecisionTree.num_obj = [2, 3]
DecisionTree.pruning = [True, False]
KNN.k = [3, 4]
KNN.distance_measure = [1 / distance,
1 - distance]

³Thornton, Chris, et al. "Auto-WEKA: Combined selection and hyperparameter optimization of classification algorithms." Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining. 2013.

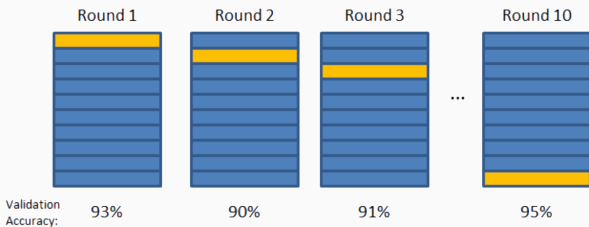
Auto-WEKA: the CASH problem

Given

- a *data-set* D divided into D_{train} , $D_{validation}$ according to k cross-validation

- $D_{train} = \{D_{train}^1, \dots, D_{train}^i, \dots, D_{train}^k\}$
- $D_{validation} = \{D_{validation}^1, \dots, D_{validation}^i, \dots, D_{validation}^k\}$
- $D_{train}^i = D \setminus D_{validation}^i$

Validation Set
Training Set



Final Accuracy = Average(Round 1, Round 2, ...)

Auto-WEKA: the CASH problem

Given

- a **data-set** D divided into $D_{train}, D_{validation}$ according to k cross-validation
 - $D_{train} = \{D_{train}^1, \dots, D_{train}^i, \dots, D_{train}^k\}$
 - $D_{validation} = \{D_{validation}^1, \dots, D_{validation}^i, \dots, D_{validation}^k\}$
 - $D_{train}^i = D \setminus D_{validation}^i$
- a set of **algorithms** $\mathcal{A} = \{A^1, \dots, A^i, \dots, A^n\}$ with associated **hyper-parameter spaces** $\{\Theta^1, \dots, \Theta^i, \dots, \Theta^n\}$

For instance:

$A^1 = \text{DecisionTree}$

$\Theta^1 = \{$
 $\text{num_obj} = [2, 3],$
 $\text{pruning} = [\text{True}, \text{False}]$
 $\}$

$A^2 = \text{KNN}$

$\Theta^2 = \{$
 $k = [3, 4],$
 $\text{distance_measure} = [1 / \text{distance}, 1 - \text{distance}]$
 $\}$

Auto-WEKA: the CASH problem

Given

- a **data-set** D divided into $D_{train}, D_{validation}$ according to k cross-validation
 - $D_{train} = \{D_{train}^1, \dots, D_{train}^i, \dots, D_{train}^k\}$
 - $D_{validation} = \{D_{validation}^1, \dots, D_{validation}^i, \dots, D_{validation}^k\}$
 - $D_{train}^i = D \setminus D_{validation}^i$
- a set of **algorithms** $\mathcal{A} = \{A^1, \dots, A^i, \dots, A^n\}$ with associated **hyper-parameter spaces** $\{\Theta^1, \dots, \Theta^i, \dots, \Theta^n\}$
- an **evaluation metric** $\mathcal{M}(A_{\theta}^i, D_{train}^i, D_{validation}^i)$

For instance:

- Accuracy
- Precision
- Recall

Auto-WEKA: the CASH problem

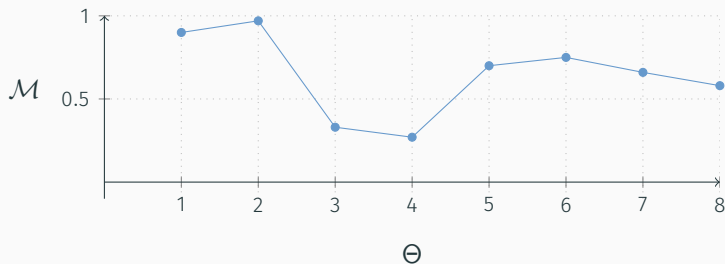
Given

- a **data-set** D divided into $D_{\text{train}}, D_{\text{validation}}$ according to k cross-validation
 - $D_{\text{train}} = \{D_{\text{train}}^1, \dots, D_{\text{train}}^i, \dots, D_{\text{train}}^k\}$
 - $D_{\text{validation}} = \{D_{\text{validation}}^1, \dots, D_{\text{validation}}^i, \dots, D_{\text{validation}}^k\}$
 - $D_{\text{train}}^i = D \setminus D_{\text{validation}}^i$
- a set of **algorithms** $\mathcal{A} = \{A^1, \dots, A^n\}$ with associated **hyper-parameter spaces** $\Theta^1, \dots, \Theta^n$
- an **evaluation metric** $\mathcal{M}(A^j, D_{\text{train}}^i, D_{\text{validation}}^i)$

We are searching for

$$A_{\theta^*}^* \in \arg \max_{A^j \in \mathcal{A}, \theta \in \Theta^j} \frac{1}{k} \sum_{i=1}^k \mathcal{M}(A_{\theta}^j, D_{\text{train}}^i, D_{\text{validation}}^i) \quad (\text{CASH})$$

Auto-Weka: CASH reformulation



| θ | Algorithm | num_obj | pruning | k | distance_measure |
|----------|--------------|---------|---------|---|------------------|
| 1 | DecisionTree | 2 | True | | |
| 2 | DecisionTree | 2 | False | | |
| 3 | DecisionTree | 3 | True | | |
| 4 | DecisionTree | 3 | False | | |
| 5 | KNN | | | 3 | 1/distance |
| 6 | KNN | | | 3 | 1-distance |
| 7 | KNN | | | 4 | 1/distance |
| 8 | KNN | | | 4 | 1-distance |

Auto-Weka: search space

| Classifier | Categorical | Numeric |
|--------------------------------|-------------|---------|
| BAYES NET | 2 | 0 |
| NAIVE BAYES | 2 | 0 |
| NAIVE BAYES MULTINOMIAL | 0 | 0 |
| GAUSSIAN PROCESS | 3 | 6 |
| LINEAR REGRESSION | 2 | 1 |
| LOGISTIC REGRESSION | 0 | 1 |
| SINGLE-LAYER PERCEPTION | 5 | 2 |
| STOCHASTIC GRADIENT DESCENT | 3 | 2 |
| SVM | 4 | 6 |
| SIMPLE LINEAR REGRESSION | 0 | 0 |
| SIMPLE LOGISTIC REGRESSION | 2 | 1 |
| VOTED PERCEPTION | 1 | 2 |
| KNN | 4 | 1 |
| K-STAR | 2 | 1 |
| DECISION TABLE | 4 | 0 |
| RFFLIB | 3 | 1 |
| M5 RULES | 3 | 1 |
| L-R | 0 | 1 |
| PART | 2 | 2 |
| D-R | 0 | 0 |
| DECISION STUMP | 0 | 0 |
| C4.5 DECISION TREE | 6 | 2 |
| LOGISTIC MODEL TREE | 5 | 2 |
| M5 TREE | 3 | 1 |
| RANDOM FOREST | 2 | 3 |
| RANDOM TREE | 4 | 4 |
| REP TREE | 2 | 3 |
| LOCALLY WEIGHTED LEARNING* | 3 | 0 |
| ADABOOST.M1* | 2 | 2 |
| ADDITIVE REGRESSION* | 1 | 2 |
| ATTRIBUTE SELECTED* | 2 | 0 |
| BAGGING* | 1 | 2 |
| CLASSIFICATION VIA REGRESSION* | 0 | 0 |
| LOGITBOOST* | 4 | 4 |
| MULTICLASS CLASSIFIER* | 3 | 0 |
| RANDOM COMMITTEE* | 0 | 1 |
| RANDOM SUBSPACE* | 0 | 2 |
| VOTING† | 1 | 0 |
| STACKING† | 0 | 0 |

Explore all the configurations is
unfeasible (786 hyper-parameters)

⇒ explore few of them but in a smart way

The table represents the considered classifiers in Auto-WEKA. Categorical and Numeric refer to the number of hyper-parameters of each kind for each classifier.

CASH resolution approaches⁴

- Model free methods
 - Grid search
 - Random search
 - Heuristics
 - Ant colony optimization
 - Particle Swarm Optimization
 - Simulate Annealing
 - Genetic algorithms
 - Multi-resolution optimization
 - Successive Halving
 - Hyper-Band
- Bayesian optimization

⁴Elshawi, R., Maher, M., Sakr, S. (2019). Automated machine learning: State-of-the-art and open challenges.

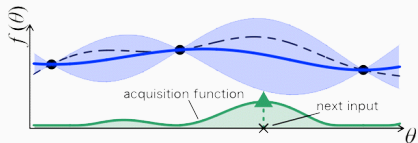
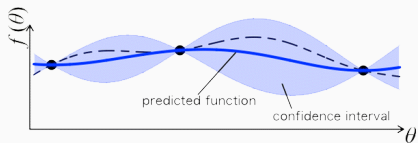
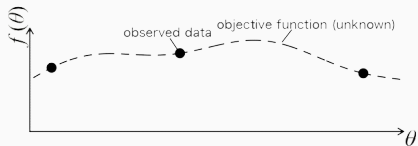
Explore all the configurations is **unfeasible**

⇒ explore few of them but in a smart way

We want to:

- divide the exploration in **iterations**
- keep track of **past evaluation scores**
- build/update a **probabilistic model**
- find promising configurations to explore

Bayesian Optimization⁵



- **objective function**: the function we want to maximize
- **observed data**: the tested hyper-parameters configurations

The **probabilistic model** consists of:

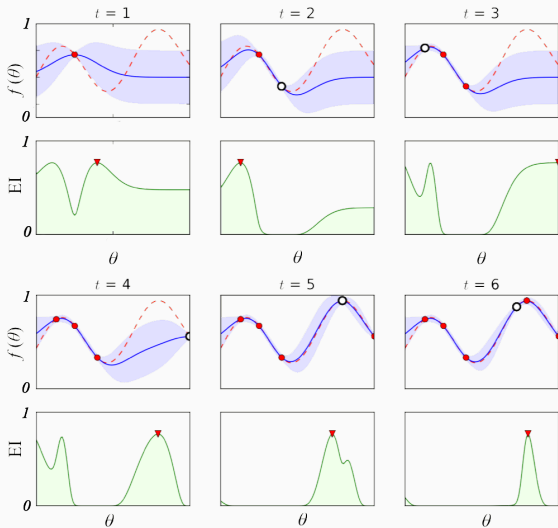
- **predicted function**, an estimation of the objective
- **confidence interval**, which indicates the possible variance

The **acquisition functions** suggests the next configuration to visit. It regulates:

- **exploitation**
- **exploration**

⁵Brochu, Eric, Vlad M. Cora, and Nando De Freitas. "A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning." (2010).

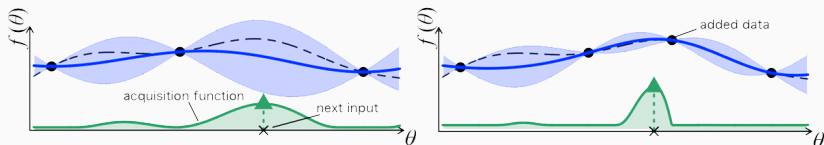
Bayesian Optimization: working example



Bayesian Optimization: SMBO

Sequential Model-Based Optimization (SMBO) is a formalization of Bayesian Optimization:

1. Evaluate some random hyper-parameters configurations
2. Build a probabilistic model
3. Exploit the model and the acquisition function to find the next hyper-parameters configuration to evaluate
4. Evaluate the hyper-parameters configuration
5. Update the probabilistic model incorporating the new results
6. Repeat steps 3–5 until the budget exceeded



The **implementations of SMBO** differ in how they construct the **probabilistic model**

- using Gaussian Process (GP)
- using Tree Parzen Estimators (TPE)
- using **Random Forest (SMAC)**⁶

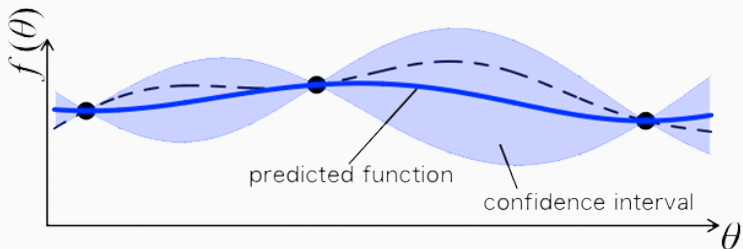
⁶F. Hutter, H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algorithm configuration. Proc. of LION-5, pages 507–523, 2011.

Bayesian Optimization: SMAC

Random Forest is not usually treated as probabilistic models.

SMAC obtains:

- the **predicted function**, as the **mean** over the predictions of its individual trees for θ
- the **confidence interval**, as the **variance** over the predictions of its individual trees for θ

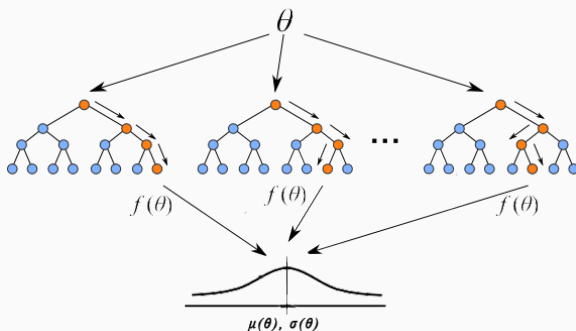


Bayesian Optimization: SMAC

Random Forest is not usually treated as probabilistic models.

SMAC obtains:

- the **predicted function**, as the **mean** over the predictions of its individual trees for θ
- the **confidence interval**, as the **variance** over the predictions of its individual trees for θ

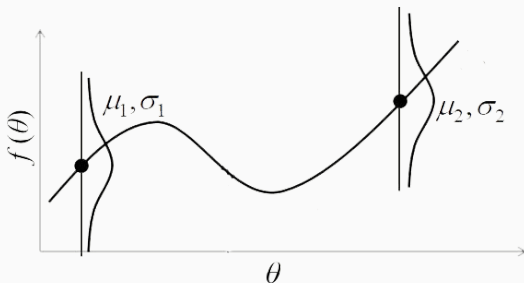


Bayesian Optimization: SMAC

Random Forest is not usually treated as probabilistic models.

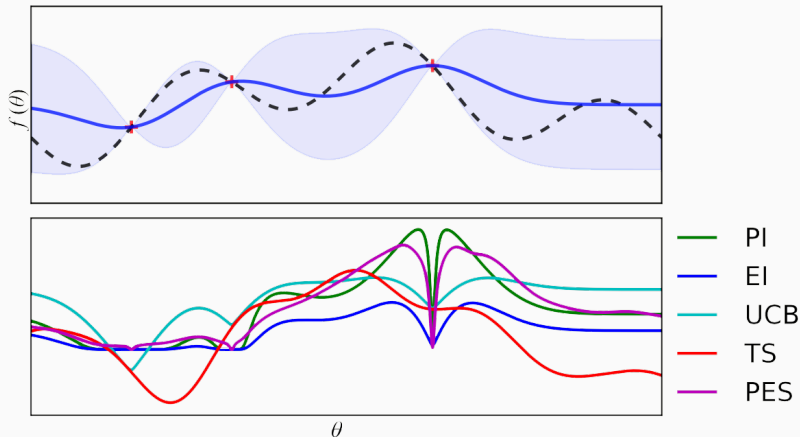
SMAC obtains:

- the **predicted function**, as the **mean** over the predictions of its individual trees for θ
- the **confidence interval**, as the **variance** over the predictions of its individual trees for θ



Bayesian Optimization: acquisition functions

The **acquisition function** is the criteria by which the next set of hyper-parameters are chosen from the surrogate function



Bayesian Optimization: Sum up

Pros:

- converge with a **low budget**
- provide **fine-grained information**

Cons:

- **slow to start** for large hyper-parameter spaces
⇒ a.k.a **cold-start problem**
- there is no optimization to reduce the **evaluation costs**

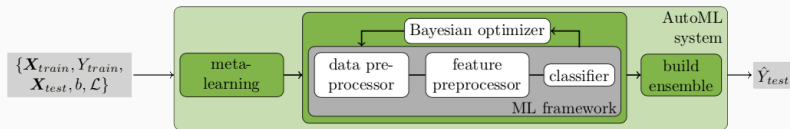
State of the art

There are three main kinds of framework:

- Cloud-Based
 - Google AutoML
 - Amazon AutoML
 - Azure AutoML
 - Data Iku
 - Data Robot
- Distributed
 - MLBase
 - TrasmogrifAI
 - MLBox
 - ATM
 - Rafiki
- Centralised
 - Auto-Weka
 - Auto-MEKA
 - Auto-Sklearn
 - HyperOpt
 - HyperOpt-Sklearn
 - TPOT
 - SmartML
 - H2O

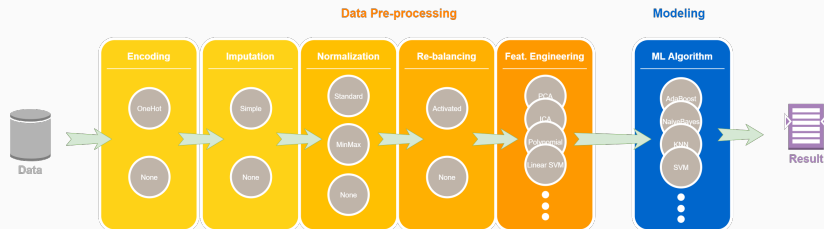
Architecture:

- Meta-learning
- Optimization
 - Scikit-learn as ML framework
 - SMAC as Bayesian optimizer
- Ensembling



⁷Feurer, Matthias, et al. "Auto-sklearn: efficient and robust automated machine learning." Automated Machine Learning. Springer, Cham, 2019. 113-134.

Auto-Sklearn: Optimization



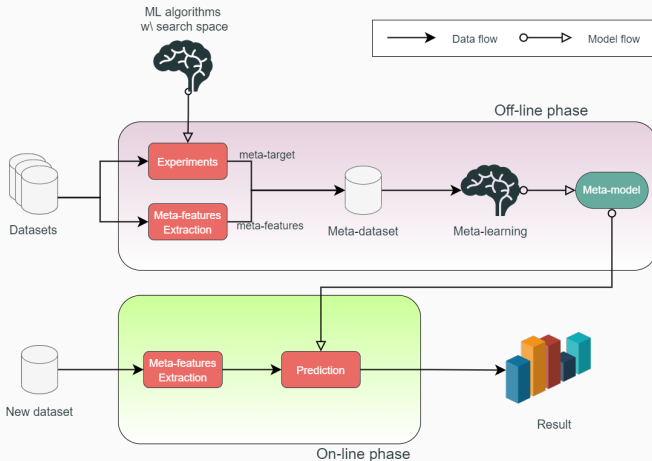
| name | # λ | cat (cond) | cont (cond) |
|-------------------------|-------------|------------|-------------|
| AdaBoost (AB) | 4 | 1 (-) | 3 (-) |
| Bernoulli naive Bayes | 2 | 1 (-) | 1 (-) |
| decision tree (DT) | 4 | 1 (-) | 3 (-) |
| extreml. rand. trees | 5 | 2 (-) | 3 (-) |
| Gaussian naive Bayes | - | - | - |
| gradient boosting (GB) | 6 | - | 6 (-) |
| kNN | 3 | 2 (-) | 1 (-) |
| LDA | 4 | 1 (-) | 3 (1) |
| linear SVM | 4 | 2 (-) | 2 (-) |
| kernel SVM | 7 | 2 (-) | 5 (2) |
| multinomial naive Bayes | 2 | 1 (-) | 1 (-) |
| passive aggressive | 3 | 1 (-) | 2 (-) |
| QDA | 2 | - | 2 (-) |
| random forest (RF) | 5 | 2 (-) | 3 (-) |
| Linear Class. (SGD) | 10 | 4 (-) | 6 (3) |

(a) classification algorithms

| name | # λ | cat (cond) | cont (cond) |
|-----------------------------|-------------|------------|-------------|
| extreml. rand. trees prepr. | 5 | 2 (-) | 3 (-) |
| fast ICA | 4 | 3 (-) | 1 (1) |
| feature agglomeration | 4 | 3 (1) | 1 (-) |
| kernel PCA | 5 | 1 (-) | 4 (3) |
| rand. kitchen sinks | 2 | - | 2 (-) |
| linear SVM prepr. | 3 | 1 (-) | 2 (-) |
| no preprocessing | - | - | - |
| nystroem sampler | 5 | 1 (-) | 4 (3) |
| PCA | 2 | 1 (-) | 1 (-) |
| polynomial | 3 | 2 (-) | 1 (-) |
| random trees embed. | 4 | - | 4 (-) |
| select percentile | 2 | 1 (-) | 1 (-) |
| select rates | 3 | 2 (-) | 1 (-) |
| one-hot encoding | 2 | 1 (-) | 1 (1) |
| imputation | 1 | 1 (-) | - |
| balancing | 1 | 1 (-) | - |
| rescaling | 1 | 1 (-) | - |

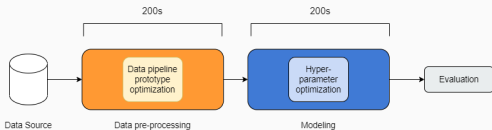
(b) preprocessing methods

Auto-Sklearn: Meta-learning



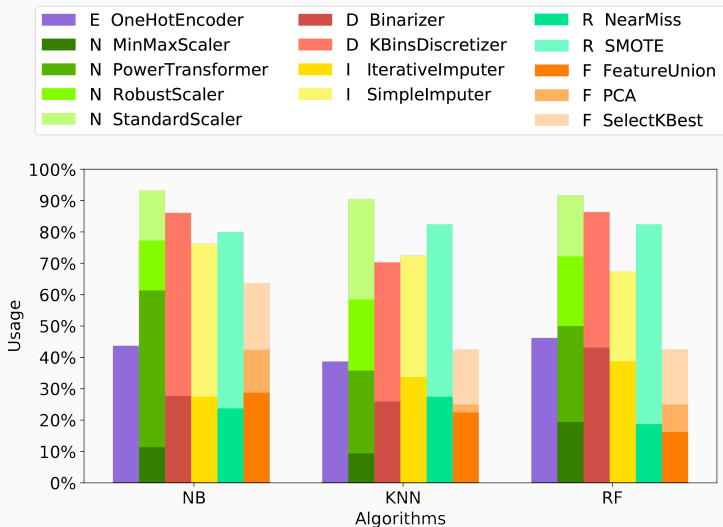
Getting insight with meta-learning

| ID | Pipeline prototype | ID | Pipeline prototype |
|----|---|----|---|
| 1 | $I \rightarrow E \rightarrow N \rightarrow D \rightarrow F \rightarrow R$ | 13 | $I \rightarrow E \rightarrow F \rightarrow N \rightarrow D \rightarrow R$ |
| 2 | $I \rightarrow E \rightarrow N \rightarrow D \rightarrow R \rightarrow F$ | 14 | $I \rightarrow E \rightarrow F \rightarrow N \rightarrow R \rightarrow D$ |
| 3 | $I \rightarrow E \rightarrow N \rightarrow F \rightarrow D \rightarrow R$ | 15 | $I \rightarrow E \rightarrow F \rightarrow D \rightarrow R \rightarrow N$ |
| 4 | $I \rightarrow E \rightarrow N \rightarrow F \rightarrow R \rightarrow D$ | 16 | $I \rightarrow E \rightarrow F \rightarrow D \rightarrow R \rightarrow N$ |
| 5 | $I \rightarrow E \rightarrow N \rightarrow R \rightarrow D \rightarrow F$ | 17 | $I \rightarrow E \rightarrow F \rightarrow R \rightarrow N \rightarrow D$ |
| 6 | $I \rightarrow E \rightarrow N \rightarrow R \rightarrow F \rightarrow D$ | 18 | $I \rightarrow E \rightarrow F \rightarrow R \rightarrow D \rightarrow N$ |
| 7 | $I \rightarrow E \rightarrow D \rightarrow N \rightarrow F \rightarrow R$ | 19 | $I \rightarrow E \rightarrow R \rightarrow N \rightarrow D \rightarrow F$ |
| 8 | $I \rightarrow E \rightarrow D \rightarrow N \rightarrow R \rightarrow F$ | 20 | $I \rightarrow E \rightarrow R \rightarrow N \rightarrow F \rightarrow D$ |
| 9 | $I \rightarrow E \rightarrow D \rightarrow F \rightarrow N \rightarrow R$ | 21 | $I \rightarrow E \rightarrow R \rightarrow D \rightarrow N \rightarrow F$ |
| 10 | $I \rightarrow E \rightarrow D \rightarrow F \rightarrow R \rightarrow N$ | 23 | $I \rightarrow E \rightarrow R \rightarrow D \rightarrow F \rightarrow N$ |
| 11 | $I \rightarrow E \rightarrow D \rightarrow R \rightarrow N \rightarrow F$ | 23 | $I \rightarrow E \rightarrow R \rightarrow F \rightarrow N \rightarrow D$ |
| 12 | $I \rightarrow E \rightarrow D \rightarrow R \rightarrow F \rightarrow N$ | 24 | $I \rightarrow E \rightarrow R \rightarrow F \rightarrow D \rightarrow N$ |



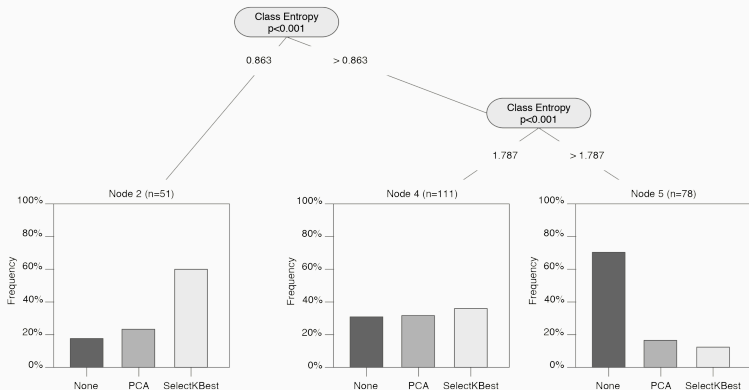
Getting insight with meta-learning

Percentage of use of transformations' operators:



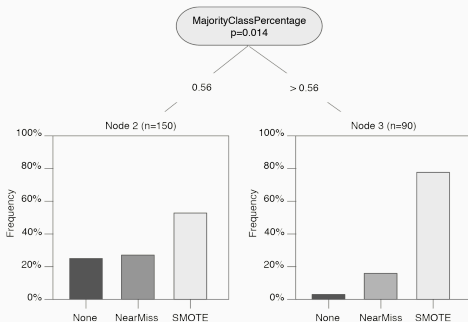
Getting insight with meta-learning

Conditional Inference Tree built for Features Engineering:



Getting insight with meta-learning

Conditional Inference Tree built for **Rebalancing**:



Meta-learning as a warm-starting procedure

Pros:

- converge with a **low budget**
- provide **fine-grained information**

Cons:

- **slow to start** for large hyper-parameter spaces
- there is no optimization to reduce the **evaluation costs**

Meta-learning as a warm-starting procedure

Pros:

- converge with a **low budget**
- provide **fine-grained information**

Cons:

- **slow to start** for large hyper-parameter spaces
⇒ a.k.a **cold-start problem**
- there is no optimization to reduce the **evaluation costs**

Meta-learning as a warm-starting procedure

Pros:

- converge with a **low budget**
- provide **fine-grained information**

Cons:

- **slow to start** for large hyper-parameter spaces
 - ⇒ a.k.a cold-start problem
- there is no optimization to reduce the **evaluation costs**
 - ⇒ **multi-fidelity optimization**

Multi-fidelity optimization

Pros:

- **Evaluate** configurations **incrementally** (e.g., folds by folds)
- **Discard non-performing** configurations

Cons:

- **Model-free** approaches

Main methods:

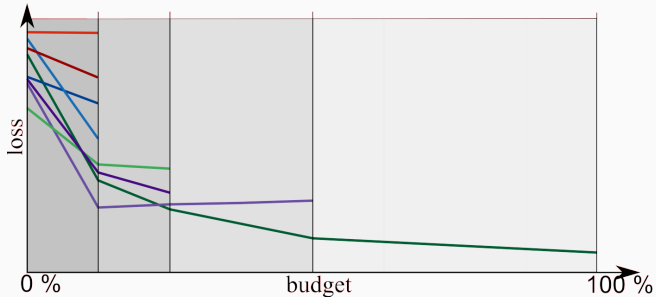
1. **Successive halving**
2. **Hyper Band**

Successive halving

Given:

- N different configurations
- a precise budget β

The evaluation starts for **all the N configurations concurrently**



At each **cut** just the **best half** of the configurations are kept

Hyper Band

Successive Halving issues:

- How we decide N number of configurations?
- How we decide the number of cuts?



Hyper Band performs frequently Successive Halving varying:

- the number of tested configurations
- the budget

Bayesian Optimization Hyper Band (BOHB)

Bayesian Optimization:

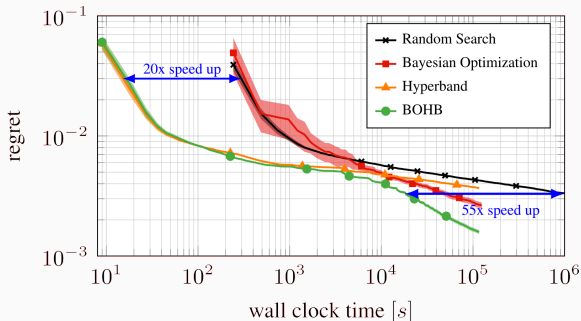
- Model-based

- really slow

Hyper Band:

- Model-free

- really fast

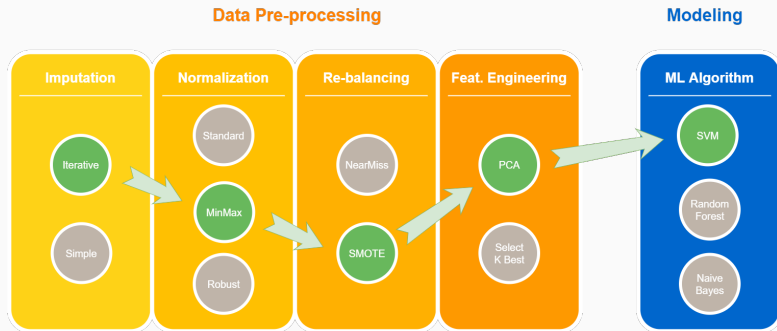


BOHB makes the most out of Bayesian Optimization and Hyper Band:

- Bayesian Optimization to **not go blindly**
- Hyper Band to **evaluate N iterations concurrently**

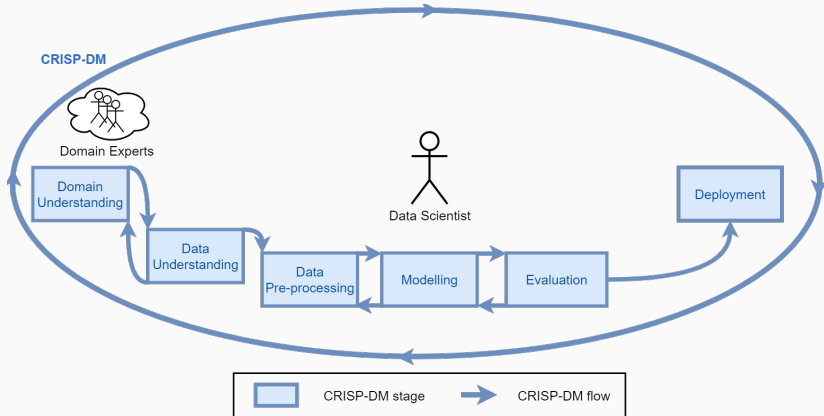
Human-centered AutoML

AutoML aims at find the best ML pipeline



- At each **step**, a **technique** must be selected
- For each **technique**, a set of **hyper-parameters** must be set
- Each **hyper-parameter** has its own **search space**

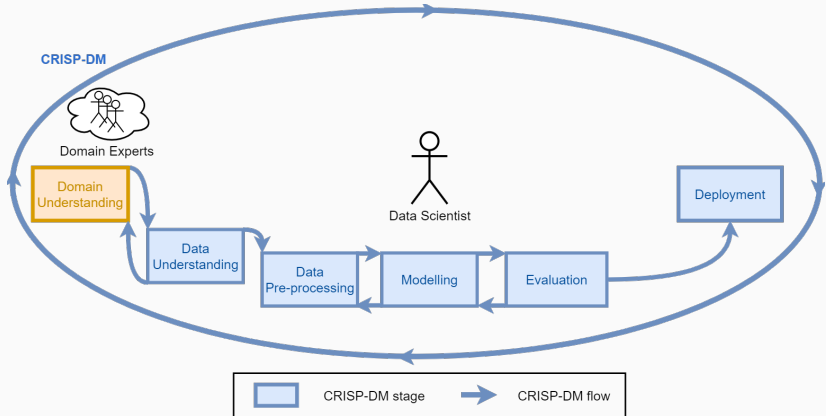
CRISP-DM: Cross Industry Standard Process for Data Mining



CRISP-DM enables the exploration of **ML Constraints**:

- domain-related;
- data-related;
- transformation-related;
- algorithm-related;
- cross-cutting (e.g., ethical, legal).

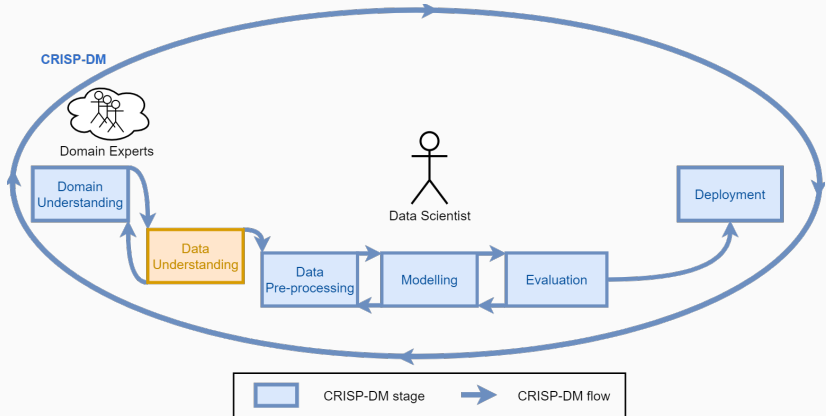
CRISP-DM: Cross Industry Standard Process for Data Mining



CRISP-DM enables the exploration of **ML Constraints**:

- domain-related;
- transformation-related;
- cross-cutting (e.g., ethical, legal).
- data-related;
- algorithm-related;

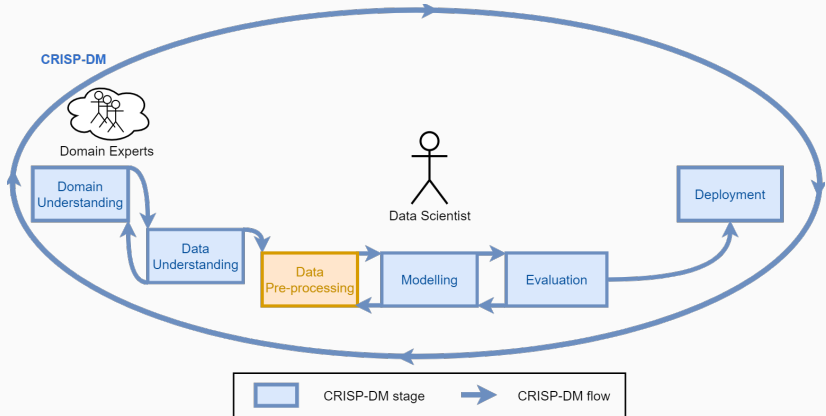
CRISP-DM: Cross Industry Standard Process for Data Mining



CRISP-DM enables the exploration of **ML Constraints**:

- domain-related;
- transformation-related;
- cross-cutting (e.g., ethical, legal).
- data-related;
- algorithm-related;

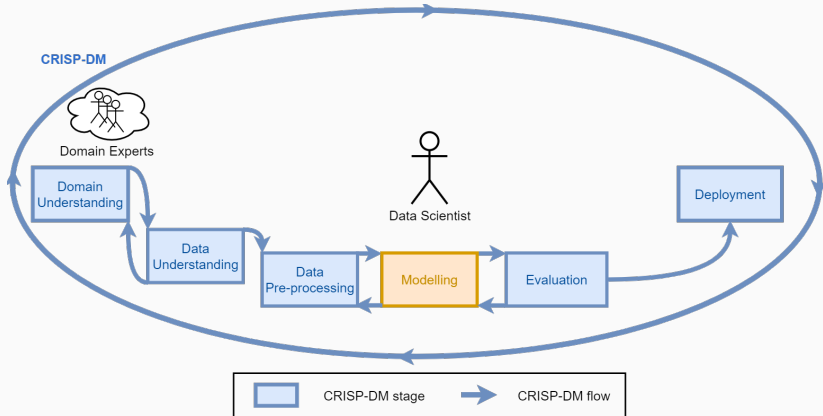
CRISP-DM: Cross Industry Standard Process for Data Mining



CRISP-DM enables the exploration of **ML Constraints**:

- domain-related;
- **transformation-related**;
- cross-cutting (e.g., ethical, legal).
- data-related;
- algorithm-related;

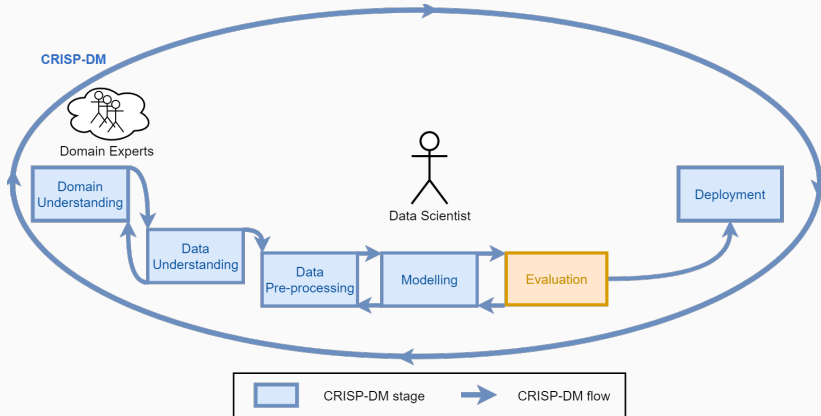
CRISP-DM: Cross Industry Standard Process for Data Mining



CRISP-DM enables the exploration of **ML Constraints**:

- domain-related;
- data-related;
- transformation-related;
- **algorithm-related**;
- cross-cutting (e.g., ethical, legal).

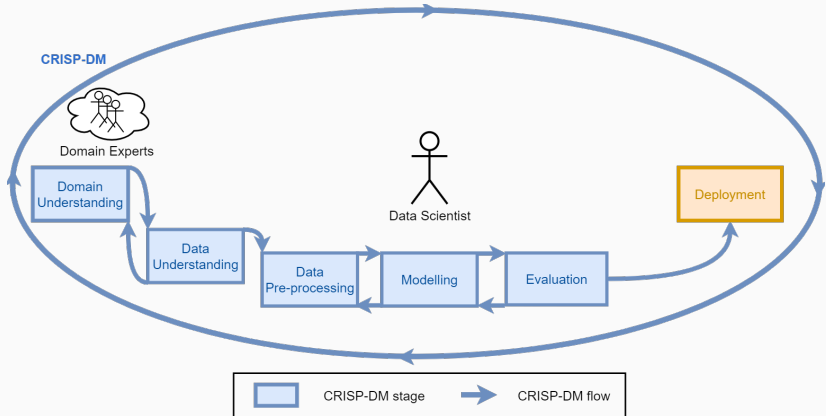
CRISP-DM: Cross Industry Standard Process for Data Mining



CRISP-DM enables the exploration of **ML Constraints**:

- domain-related;
- data-related;
- transformation-related;
- algorithm-related;
- cross-cutting (e.g., ethical, legal).

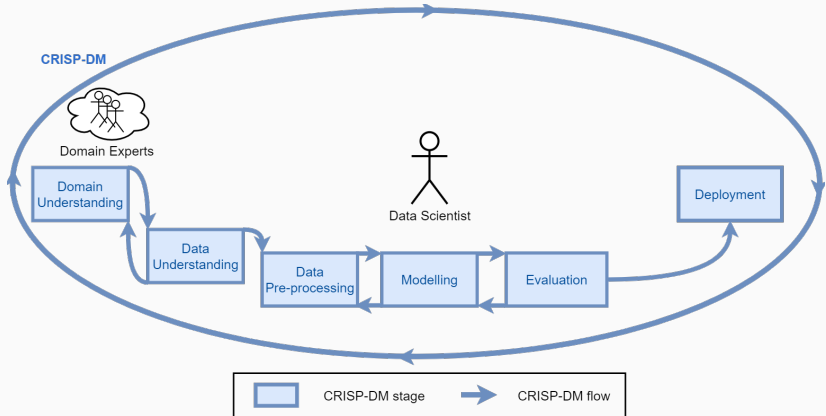
CRISP-DM: Cross Industry Standard Process for Data Mining



CRISP-DM enables the exploration of **ML Constraints**:

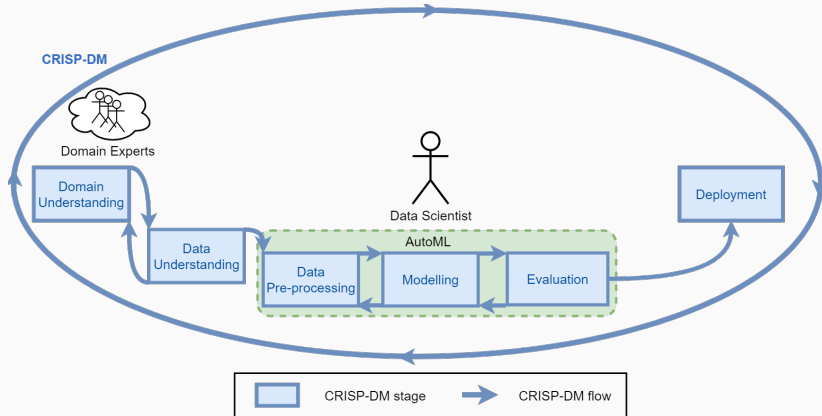
- domain-related;
- data-related;
- transformation-related;
- algorithm-related;
- cross-cutting (e.g., ethical, legal).

CRISP-DM: Cross Industry Standard Process for Data Mining



CRISP-DM enables the exploration of **ML Constraints**:

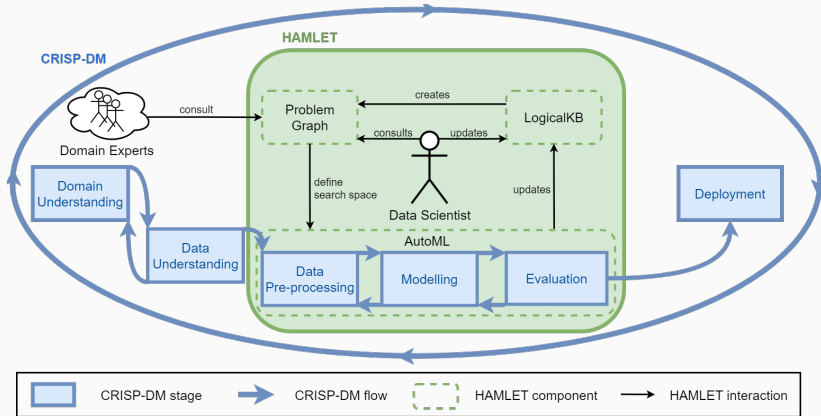
- domain-related;
- data-related;
- transformation-related;
- algorithm-related;



AutoML aims at automating the ML pipeline instantiation:

- it is difficult to consider all the constraints together;
- it is not transparent;
- it doesn't allow a proper knowledge augmentation.

HAMLET: Human-centric AutoML via Logic and Argumentation



HAMLET leverages :

- **Logic** to give a structure to the knowledge;
- **Argumentation** to deal with inconsistencies, and revise the results.

Questions?