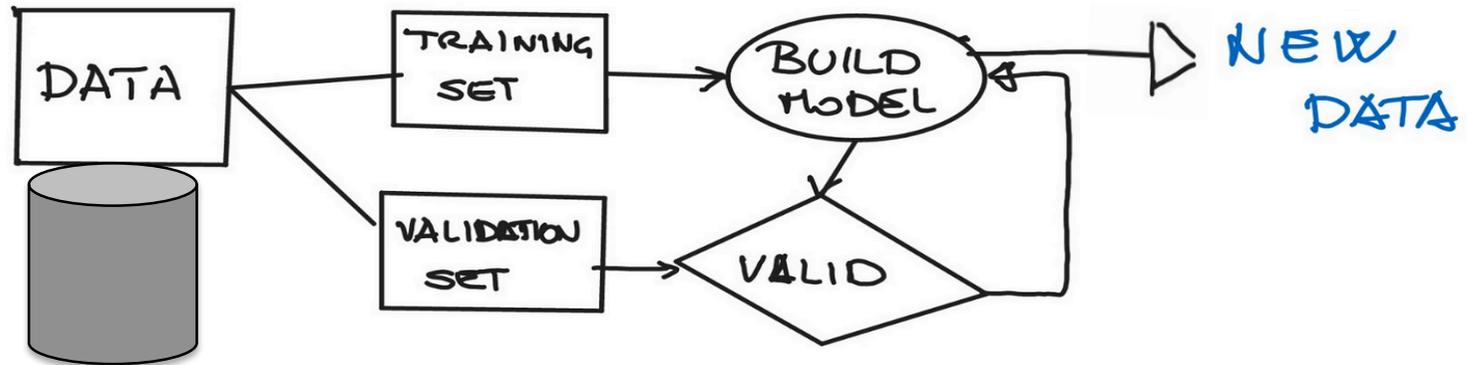# Data Curation for ML:
# Toward a Principled Approach

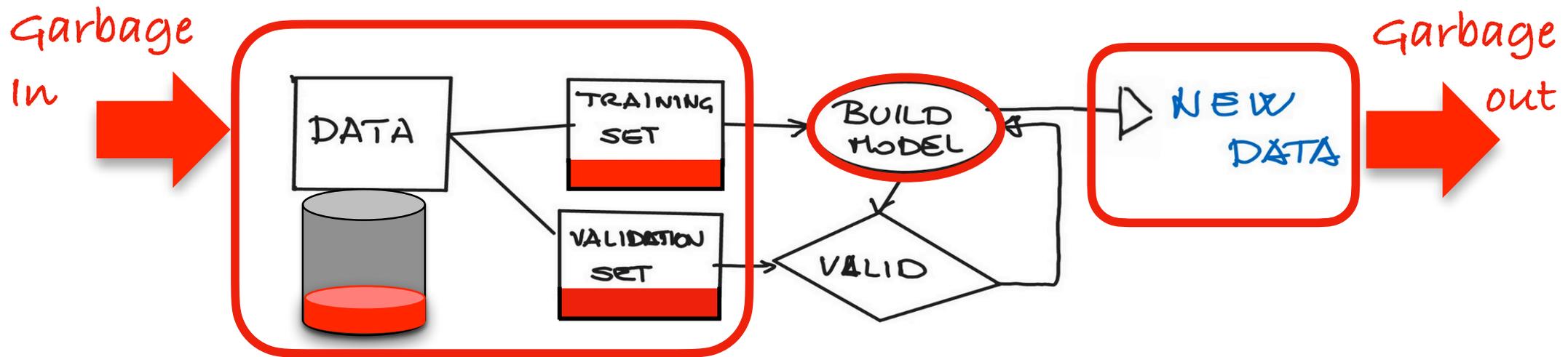**Laure Berti-Equille**

laure.berti@ird.fr

Espace-Dev, IRD, Univ Montpellier, Univ Guyane, Univ La Réunion, Univ Antilles, Univ Nouvelle Calédonie, Montpellier France
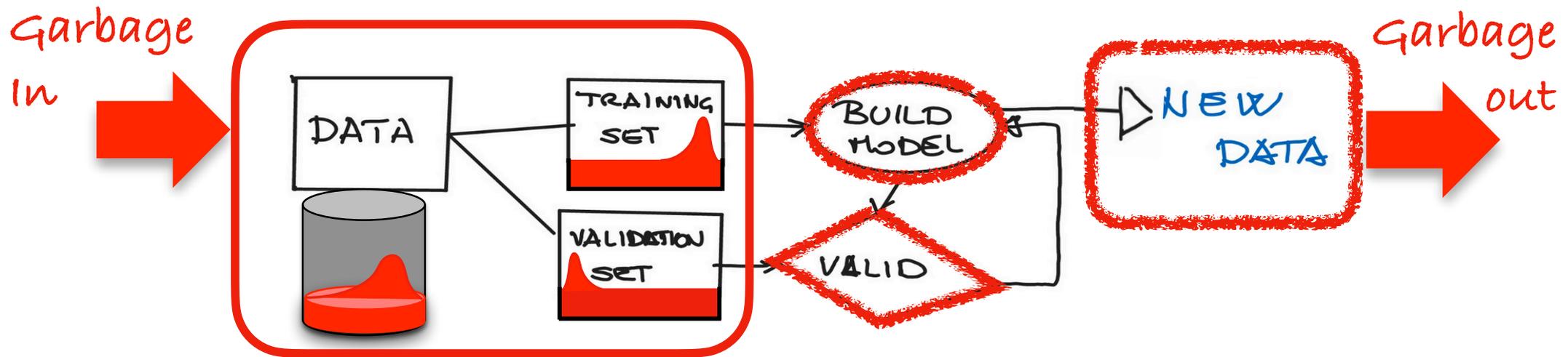
# Learning from dirty data is risky

# Learning from dirty data is risky



Garbage in

DATA

TRAINING SET

VALIDATION SET

BUILD MODEL

VALID

NEW DATA

Garbage out

# Learning from dirty data is risky



Glitch types and distributions can be very different in the datasets used for training, testing, and validation and they affect accuracy of ML models in different ways.

# Two complementary approaches

**INTERVENE**

**How to efficiently fix the data:**
✦ Detect the anomalies
✦ Correct them with minimal cost (domain expert intervention, time, external master data, etc.)
✦ Select the repair/preparation strategies that will maximize the ML result quality

**MITIGATE**

**How to reduce the impact of dirty data:**
✦ Robustify the ML algorithms and apply ML ensembling strategies
✦ Use AutoML to find optimal parameter setting
✦ Select portions of the data and/or augment the data

# Two complementary approaches

## INTERVENE

**How to efficiently fix the data:**
- ✦ Detect the anomalies
- ✦ Correct them with minimal cost (domain expert intervention, time, external master data, etc.)
- ✦ Select the repair/preparation strategies that will maximize the ML result quality

## MITIGATE

**How to reduce the impact of dirty data:**
- ✦ Robustify the ML algorithms and apply ML ensembling strategies
- ✦ Use AutoML to find optimal parameter setting
- ✦ Select portions of the data and/or augment the data

# Outline

1. **Detection of data quality problems**
   Profiling data quality

2. **Data cleaning**
   Leveraging the patterns of glitches

3. **Data preparation strategies for ML**
   Learning to clean and prepare the data

# Outline

1. **Detection of data quality problems**
   Profiling data quality

2. **Data cleaning**
   Leveraging the patterns of glitches

3. **Data preparation strategies**
   Learning to clean and prepare the data

# Data Quality Problems



**DATA TYPES**

0101010101

ACACGTGT

John Doe

High
Medium
Low

**RELATIONSHIPS**

| |
| --- |
| Structural (record) |
| Sequential |
| Graph-based |
| Temporal |
| Spatial |
| Spatio-Temporal |

**DATA QUALITY PROBLEMS**

| TYPE | CARDINALITY |
| --- | --- |
| Missing Data<br>Anomalous Data<br>Duplicate Data<br>Inconsistent Data<br>Obsolete Data<br>Incorrect data | Single-Point<br>Collection |

| DETECTION MODE |
| --- |
| Model-based |
| Data distribution-based |
| Constraint-based |
| Pattern-based |

# Existing approaches for detecting/fixing DQ problems

## Declarative
- Data debugging
- Checking data assertions
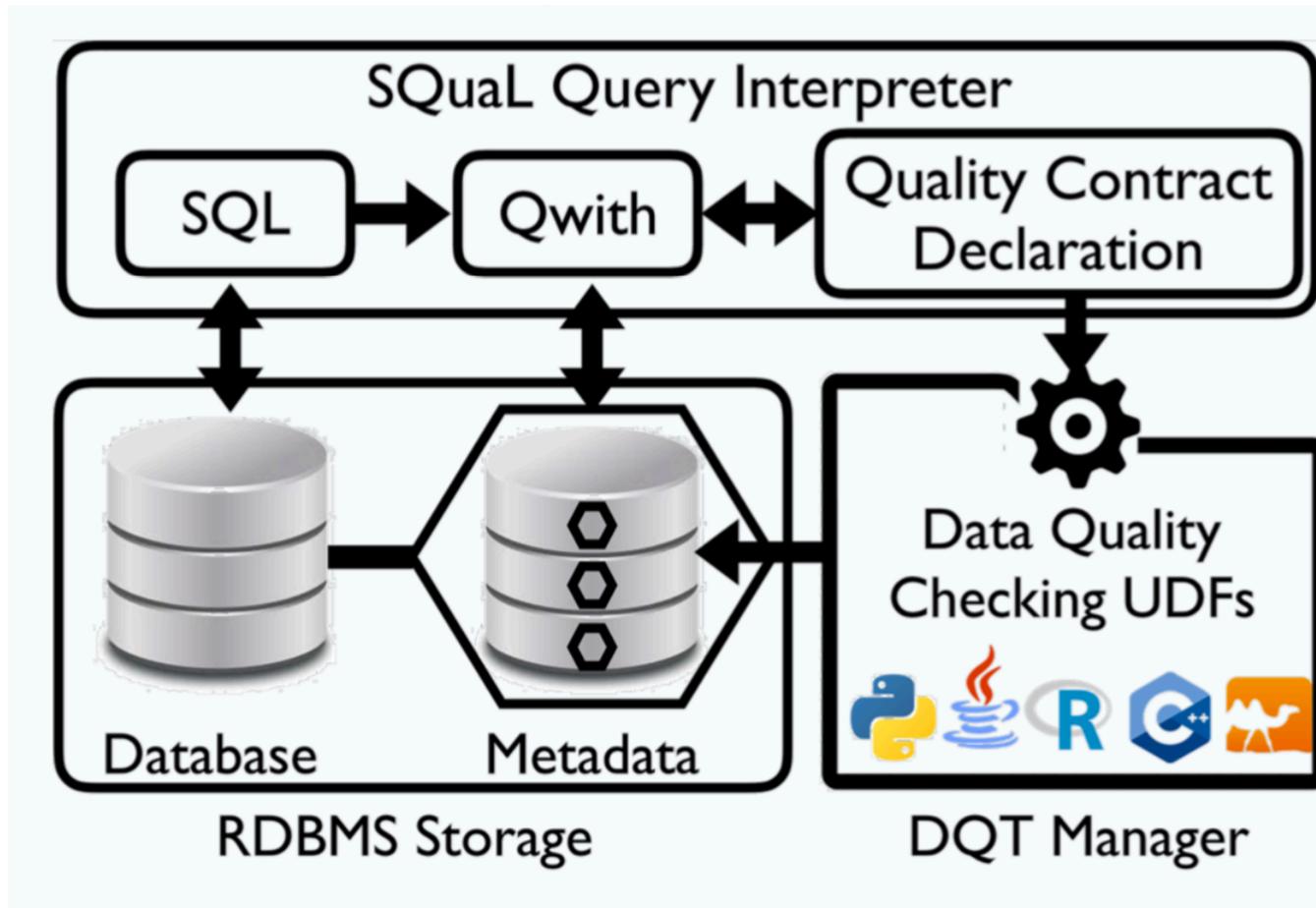- Transform

## ML-based
Learn from clean data and replace

# Declarative Approaches

**Checking data assertions and transform**

✦ **Deequ** [Schelter et al., VLDB 2018] requires cloud infrastructure and manual integration into training and serving systems; dependent on Apache Spark

✦ **TensorFlow Data Validation** (TFDV) [Caveness et al., SIGMOD 2020] integrated with Google TFX difficult to use outside of these platforms

✦ Lightweight Python-based approaches like **great_expectations** (https://greatexpectations.io) or **hooqu** (https://github.com/mfcabrera/hooqu) not integrated with the ML development process

# Declarative data profiling with MeSQuaL



SQuaL Query Interpreter

SQL → Qwith ↔ Quality Contract Declaration

Database — Metadata

RDBMS Storage

Data Quality Checking UDFs

DQT Manager

https://github.com/ucomignani/MeSQuaL
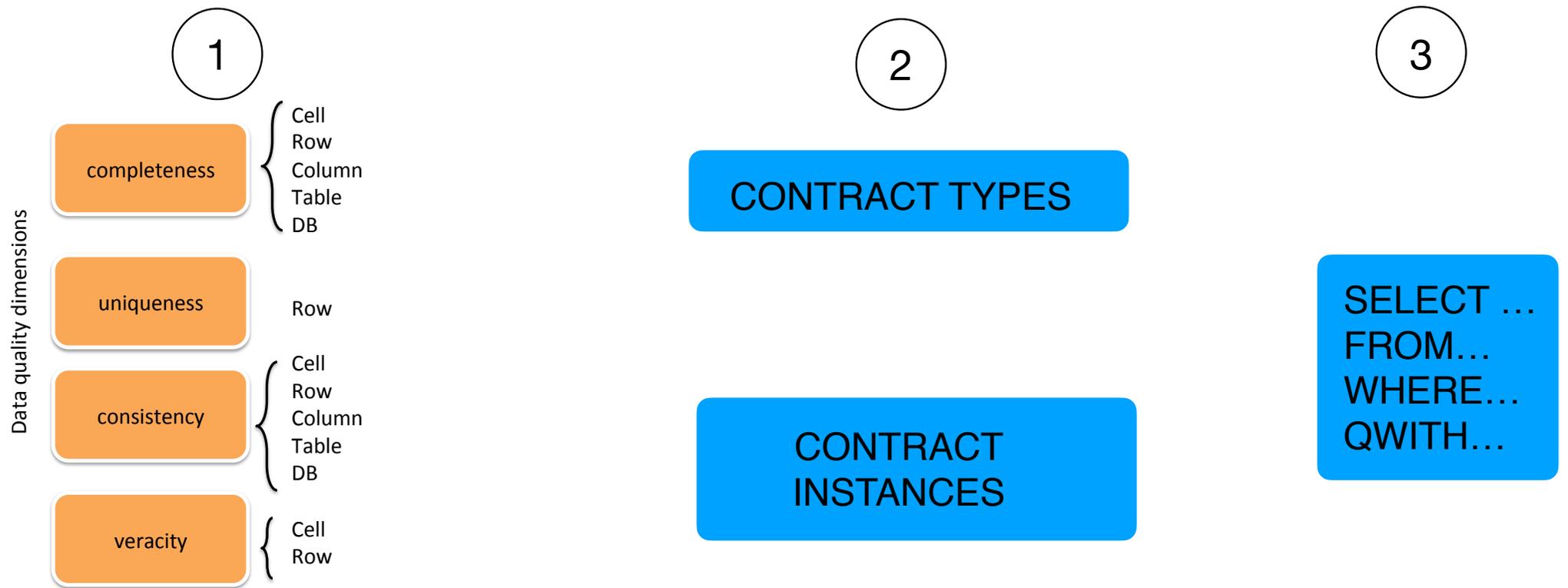
# MeSQuaL Key Concepts

*Flexible declarative data quality profiling with UDFs*



**Procedural approach with UDFs**

**Declarative approach**

**Extended query**

13

# MeSQuaL Examples

## DECLARATION

```
CREATE CONTRACTTYPE StatTests (
    autocorrelation BY FUNCTION 'durbinWatsonTest.py' LANGUAGE PYTHON,
    multicollinearity BY FUNCTION 'varInflationFactor.py' LANGUAGE PYTHON,
    heteroscedasticity BY FUNCTION 'BreuschPaganTest.py' LANGUAGE PYTHON,
    KMerrorNormality BY FUNCTION 'KolmogorovSmirnov.py' LANGUAGE PYTHON,
    SWerrorNormality BY FUNCTION 'ShapiroWilkTest.py' LANGUAGE PYTHON);
```

```
CREATE CONTRACT RegressionAssumptions (
    StatTests.autocorrelation > 0
    AND StatTests.autocorrelation < 4
    AND StatTests.multicollinearity <= 4
    AND StatTests.heteroscedasticity < 0.05
    AND StatTests.SWerrorNormality < 0.05);
```

```
CREATE CONTRACTTYPE CheckQDB (
    completeness BY FUNCTION 'completeness.py' LANGUAGE PYTHON,
    uniqueness BY FUNCTION 'uniqueness.py' LANGUAGE PYTHON,
    consistency BY FUNCTION 'consistency.py' LANGUAGE PYTHON,
    outlyingness BY FUNCTION 'outlyingness.py' LANGUAGE PYTHON);
```

```
CREATE CONTRACT CheckBeforeAnalysis (
    RegressionAssumptions
    AND CheckQDB.consistency > 0.9
    AND CheckQDB.outlyingness < 0.2);
```

## MANIPULATION

| | |
|---|---|
| AoT | `{ SELECT * FROM ChicagoDataset } QWITH CheckQDB.completeness> 0.95;` |
| | `{ SELECT * FROM ChicagoDataset } QWITH CheckBeforeAnalysis AND RegressionAssumptions;` |
| | `{ SELECT timestamp, node_id,value_raw,valuehrf FROM ChicagoDataset WHERE ChicagoDataset.sensor = 'o3' } QWITH CheckBeforeAnalysis AND CheckQDB.completeness> 0.95;` |
| MIMIC-III | `{ SELECT * FROM Admissions } QWITH CheckQDB.completeness> 0.95;` |
| | `{ SELECT * FROM Admissions WHERE Admissions.insurance = 'Private' } QWITH CheckBeforeAnalysis AND CheckQDB.completeness> 0.95;` |
| | `{ SELECT gender, dob, admittime FROM Admissions INNER JOIN (SELECT * FROM Patients WHERE dob < '2090-12-12 00:00:00' QWITH CheckQDB.completeness> 0.95) as Pat ON Admissions.subject_id=Pat.subject_id; } QWITH CheckQDB.completeness> 0.95;` |

# MeSQuaL GUI

# ML-based Approaches

**Learn from clean data and replace/repair**

- Pattern enforcement
  - Syntactic patterns (date formatting)
  - Semantic patterns (name/address)
- Value update to satisfy a set of rules, constraints, FDs, CFDs, Denial Constraints (DCs), Matching Dependencies (MDs) with minimal number of changes.
- Value replacement
- Entity resolution

**EXAMPLES**

✦ SCARE: Scalable Automatic Repair

✦ On-demand ETL [Yang et al., VLDB'15]

✦ ActiveClean [Krishnan et al., VLDB'16]

✦ HoloClean [Rekatsinas et al., VLDB 2017]

✦ Deep learning for Entity Resolution

✦ Transformers for data prep

# SCARE: SCalable Automatic Repair

[Yakout, Berti-Equille, Elmagarmid, SIGMOD 2013]

*Goal: Find the repair that would maximize the sum of the probabilities of the values co-occurrence (i.e., association strength between predicted and reliable values) under a certain update cost budget.*

**Reliable**    **Flexible**



1. Modeling Dependency and Predicting Updates

2. Data Partitioning

3. Tuple Repair Selection

Database Table

Partitioning Functions $f_1$ ... $f_i$ ... $f_n$

Machine Learning Models

Predicted updates    Predicted updates

Candidate Tuple Repairs

*Value predictions for Flexible Attributes E1, E2, E3*

# On-demand ETL with Lenses

[Yang et al. , VLDB'15]

Specification of Lens with classifiers from the massive online analysis (MOA) framework for Domain Constraint Repair (DCR).



```
CREATE LENS SaneProduct AS SELECT * FROM Product
  USING DOMAIN_REPAIR( category string NOT NULL,
                       brand    string NOT NULL );
```

| id | name | brand | category |
|---|---|---|---|
| P123 | Apple 6s, White | $Var('X', \text{R1})$ | phone |
| P124 | Apple 5s, Black | $Var('X', \text{R2})$ | phone |
| P125 | Samsung Note2 | Samsung | phone |
| P2345 | Sony 60 inches | $Var('X', \text{R4})$ | $Var('Y', \text{R4})$ |
| P34234 | Dell, Intel 4 core | Dell | laptop |
| P34235 | HP, AMD 2 core | HP | laptop |

# HoloClean

[Rekatsinas et al., VLDB 2017]

https://github.com/HoloClean/HoloClean

*HoloClean generates a factor graph capturing co-occurrences, correlations based on a set of constraints and external evidences. It uses SGD to learn parameters and infer the marginal distribution of unknown variables with Gibbs sampling.*

Each cell is a random variable

Value co-occurrences capture data statistics

Constraints introduce correlations

$c1: Zip \rightarrow City$

| | Address | City | State | Zip |
|---|---|---|---|---|
| t1 | 3465 S Morgan ST | *Chicago* | IL | *60608* |
| t2 | 3465 S Morgan ST | Chicago | IL | *60609* |
| t3 | 3465 S Morgan ST | Chicago | IL | *60609* |
| t4 | 3465 S Morgan ST | *Cicago* | IL | *60608* |

"Address= 3465 S Morgan St"

**t1.City**      **t1.Zip**

c1

**t4.City**      **t4.Zip**

◯ : Unknown (to be inferred) RV

■ : Factor (encodes correlations)

Denial constraints:

$\forall t_1, t_2 \in D : \neg(t1[Zip] = t2[Zip] \wedge t1[City] \neq t2[City])$

$\forall t_1, t_2 \in D : \neg(t1[Zip] = t2[Zip] \wedge t1[State] \neq t2[State])$

19

# BoostClean

[Krishnan et al., 2017]

*BoostClean selects an ensemble of methods (statistical and logic rules) for error detection and for repair combinations using statistical boosting.*

Test data   Training data

Boost & Clean
train model

Detector Library
IsoDetect
Repair Library

Detectors
Repairs

Test Accuracy Evaluator

$L^*$

Deployer

Robust Classifier $C^*$

**Algorithm ˉ: Boost-and-Clean Algorithm**

**Data:** (X, Y)

1. Initialize $W_i^{(1)} = \frac{1}{N}$
2. $\mathcal{L}$ generates a set of classifiers $\mathcal{C}\{C^{(0)}, C^{(1)}, ..., C^{(k)}\}$ where $C^{(0)}$ is the base classifier and $C^{(1)}, ..., C^{(k)}$ are derived from the cleaning operations.
3. **for** $t \in [1, T]$ **do**
4.    $C_t = $ Find $C_t \in \mathcal{C}$ that maximizes the weighted accuracy on the test set. $\epsilon_t = $ Calculate weighted classification error on the test set $\alpha_t = \ln(\frac{1-\epsilon_t}{\epsilon_t})$
   $W_i^{(t+1)} \propto W_i^{(t)} e^{-\alpha_t y_i C_t(x_i)}$: down-weight correct predictions, up-weight incorrectly predictions.
5. **return** $C(x) = \text{sign}(\sum_t^T \alpha_t C_t(x))$

# Record Linkage (RL): Generic Workflow

Database R

| Name | SSN | Addr |
|------|-----|------|
| Will Forth | 354-564-339 | Ada Bd |
| Jacky Khan | 435-232-129 | Marple Street |
| Dom Hack | 235-575-689 | Main Street |
| ... | ... | ... |

Database S

| Name | SSN | Addr |
|------|-----|------|
| Jack Khan | 435-223-129 | Marple St |
| Hans Ford | 354-564-339 | Clover Bd |
| Tom Hack | 235-557-689 | Main St |
| ... | ... | ... |

R X S

[Fellegi, Sunter, 1969]
[Christen, 2012]

Cleaning Standardization

Attribute selection

Blocking

- Hashing
- Sorted keys
- Sorted NN
- (Multiple) Windowing
- Clustering

{pairs}

Record pair comparison

- Token-based : N-grams…
- Distance-based: Jaro, Edit, Levenshtein, Soundex
- Domain-dependent

{comparison vectors}

Decision Model

**Linkage decision**: $RL(pair) = \dfrac{P(vector \mid pair \in Match)}{P(vector \mid pair \in Non\ Match)}$

$L$      $U$

$RL(pair)$

| Non Match | Potential Match | Match |

21

# Human-In-The Loop for Entity Matching

**Magellan project: Lessons learnt for How-to Guide for EM**

# Deep learning for Entity Resolution

Record pair → Relevant word extraction → Word embedding → DNN → Binary classification → Match / UnMatch

GloVE          LSTM-RNN

tuple t  | A1 ... Ap ... Am |

DeepER [Ebraheem et al., Arxiv 2017]

Embedding lookup layer

Composition (avg, LSTM) layer

Similarity layer

Dense layer

Classification layer

Words

w1
wi
wj

v1
vi
vk

Words

tuple t'  | A1 ... Ap ... Am |

23

# Outline

1. **Detection of data quality problems:**
   Profiling data quality

2. **Data cleaning**
   Leveraging the patterns of glitches

3. **Data preparation strategies:**
   Learning to clean and prepare the data

# SNMP Data Analysis

- Periodic inbound and outbound traffic measurements from interfaces of network devices
- 10 attributes, every 5 minutes, over 4 weeks
- Axes transformed for plotting

# SNMP Data Analysis



Outliers

Outliers

Interfaces

Utilization_Out
Utilization_In

Bytes_Out

Bytes_In

Memory

CPU
Latency
Syslog_Events

CPU_Poll

Missing

Duplicates

- Periodic inbound and outbound traffic measurements from interfaces of network devices
- 10 attributes, every 5 minutes, over 4 weeks
- Axes transformed for plotting

Time in Frames

26

# SNMP Data Analysis



**Outliers**

**Outliers**

Interfaces

Utilization_Out

- Periodic inbound and outbound traffic measurements from interfaces of network
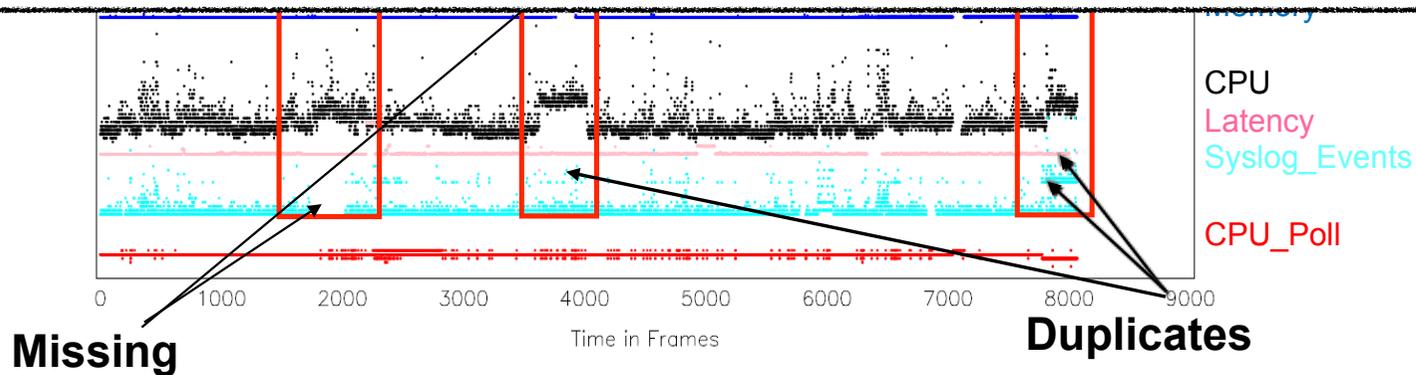
1. Detect patterns of multivariate, concomitant data anomalies

2. Use the anomaly patterns for consistent cleaning

Memory

CPU
Latency
Syslog_Events

CPU_Poll

0   1000  2000  3000  4000  5000  6000  7000  8000  9000

Time in Frames

**Missing**

**Duplicates**

27

# Understanding Complex Glitch Patterns

## *Benefits*

- A common root cause can generate correlated data errors
- In-depth anomaly analysis could help for:
  - Characterizing anomaly sources, processes, and propagation mechanisms
  - Systematizing data cleaning

## *Current methods*

- Make unrealistic assumptions (e.g., MAR)
- Treat glitches in isolation
- Are one-shot approaches (no reiteration between detection and cleaning)

  Data cleaning and preprocessing may introduce new errors and distortions.

# Detection-Exploration-Cleaning Framework

## *Input:*

**Dataset**

**Detection methods**

- ● **M**issing value
- ● **O**utlying value
- ● **I**nconsistent value
- ▶ **D**uplicate record

+

**Detection of patterns of anomalies**

# Detection-Exploration-Cleaning Framework

## *Input:*

### Dataset



**Detection methods**

- ● **M̲issing value**
- ● **O̲utlying value**
- ● **I̲nconsistent value**
- ▶ **D̲uplicate record**

\+

**Detection of patterns of anomalies**

### Specifications of the ideally preprocessed dataset



e.g., less than 5% of anomalies, should preserve the median of the original data distributions, etc.

# Detection-Exploration-Cleaning Framework

## Input:

### Dataset



**Detection methods**
- ● **M**issing value
- ◯ **O**utlying value (pink)
- ◯ **I**nconsistent value (blue)
- ▶ **D**uplicate record (orange)
+
**Detection of patterns of anomalies**

### Candidate cleaning strategies

**Cleaning**

- **Deletion**
  - **List-wise**
  - **Pair-wise**
- **Value replacement by imputation**
  - **Single**
    - **Median**
    - **EM**
    - **Logistic**
  - **Multiple**
    - **Discriminant**
    - **MCMC**
    - **Regression**
    - **Cluster mean**
- **Record replacement by fusion/selection**

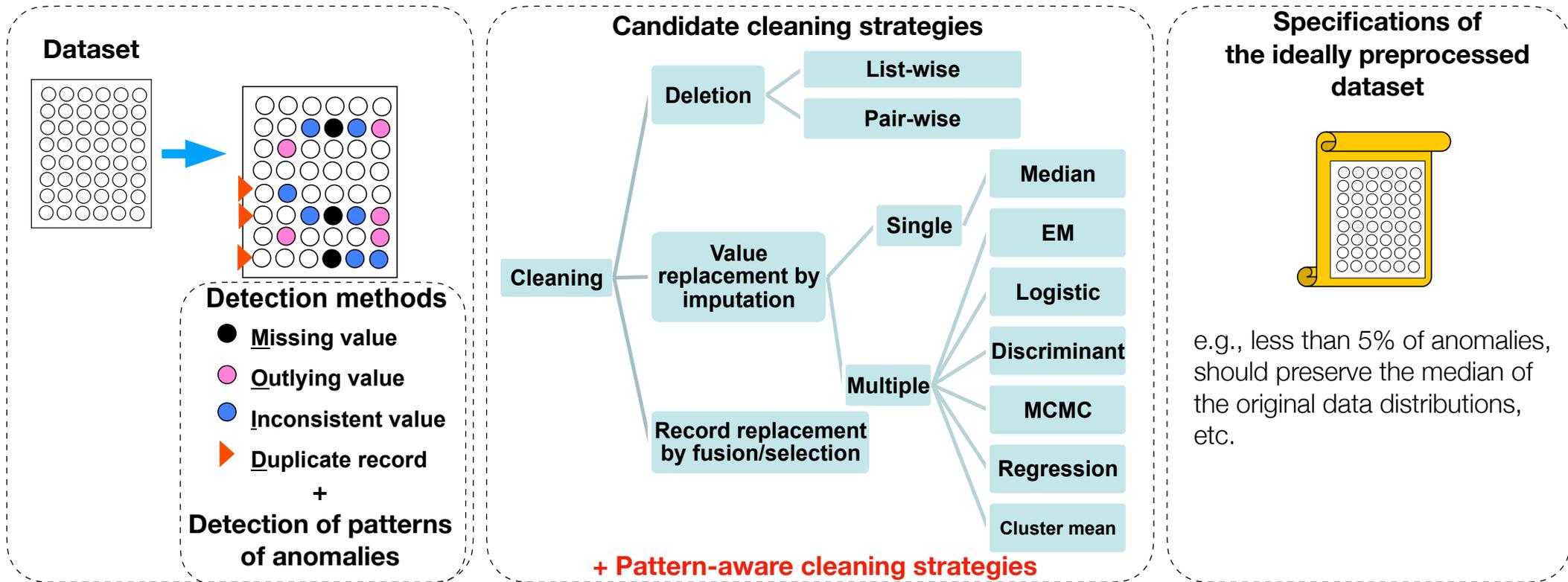**+ Pattern-aware cleaning strategies**

### Specifications of the ideally preprocessed dataset



e.g., less than 5% of anomalies, should preserve the median of the original data distributions, etc.
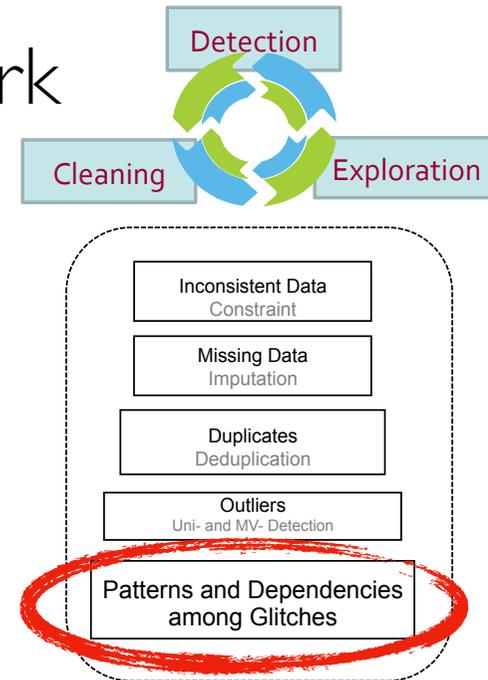
# Detection-Exploration-Cleaning Framework

[Berti-Equille, Dasu,Srivastava, ICDE 2011]

## Detection

## Cleaning    Exploration

Inconsistent Data
Constraint

Missing Data
Imputation

Duplicates
Deduplication

Outliers
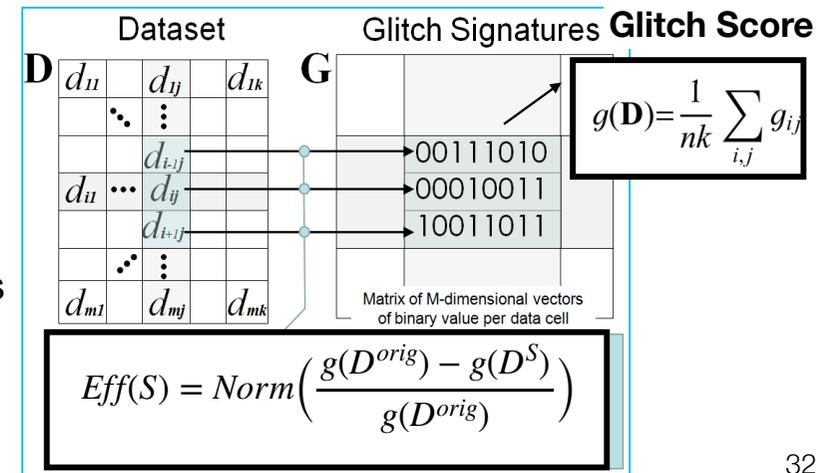Uni- and MV- Detection

Patterns and Dependencies
among Glitches

*Problem Statement:*

*Find the quantitative cleaning strategy $B$ composed of $M$ methods among the candidate strategies $S$ such that its resulting dataset $D^B$ is the closest to an ideal dataset $D*$ specified from $D$ as*

$$D^B = \arg\min_{\{s \in S\}} \left( \mathrm{dist}(D^s, D*) \right)$$

subject to $Cost(s) \leq U$ and $Eff(s) \geq \Gamma > 0$

**dist** is the Kullback-Leibler distance between two data distributions
**U**  is a pre-defined upper bound for the cost of strategy  $s$
**Γ**  is the lower bound of $Eff(s)$, the effectiveness of strategy $s$

Dataset        Glitch Signatures **Glitch Score**

| **D** | $d_{11}$ | | $d_{1j}$ | | $d_{1k}$ | **G** |
|---|---|---|---|---|---|---|
| | | | | | | |
| | | | $d_{i-1j}$ | | | → 00111010 |
| $d_{i1}$ | ... | | $d_{ij}$ | | | → 00010011 |
| | | | $d_{i+1j}$ | | | → 10011011 |
| | | | | | | |
| $d_{m1}$ | | | $d_{mj}$ | | $d_{mk}$ | |

$$g(\mathbf{D}) = \frac{1}{nk} \sum_{i,j} g_{ij}$$

Matrix of M-dimensional vectors
of binary value per data cell

$$Eff(S) = Norm\left( \frac{g(D^{orig}) - g(D^S)}{g(D^{orig})} \right)$$

32

# Experiments

## Real-world and semi-synthetic data

- **EPO Dataset:** 754,075 records, 4 non-key attributes (string, categorical and numerical data)
- **Intel Berkeley Research lab Dataset:** 2,313,682 million readings, 8 attributes (timestamp, sensorID, temperature, light, voltage) collected every 31 seconds from 54 sensors deployed in the between February 28th and April 5th, 2
- **SNMP Dataset:** (8,632 tuples, 11 variables) collected every 5 minutes during one month (timestamps, categorical and numerical values)

## Comparison of various cleaning strategies

- Cost-based
- Effectiveness-based
- Resource-driven to treat just p% of glitches (DEC-RD)
- Specification-driven to treat a particular glitch type (DEC-SD)
- Pattern-based (DEC-PD)

# Experimental results

**SNMP**



➡ **Pattern discovery always improves the accuracy when glitch percentage increases. Effectiveness is improved by + 8% in average.**

➡ *61% of the best strategies are pattern-based.*

# Outline

1. **Detection of data quality problems:**
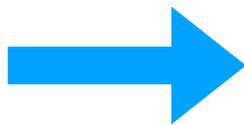      Profiling data quality with MeSQuaL

2. **Data cleaning**
      Leveraging the patterns of glitches

3. **Data preparation strategies**
      Learning to clean and prepare the data
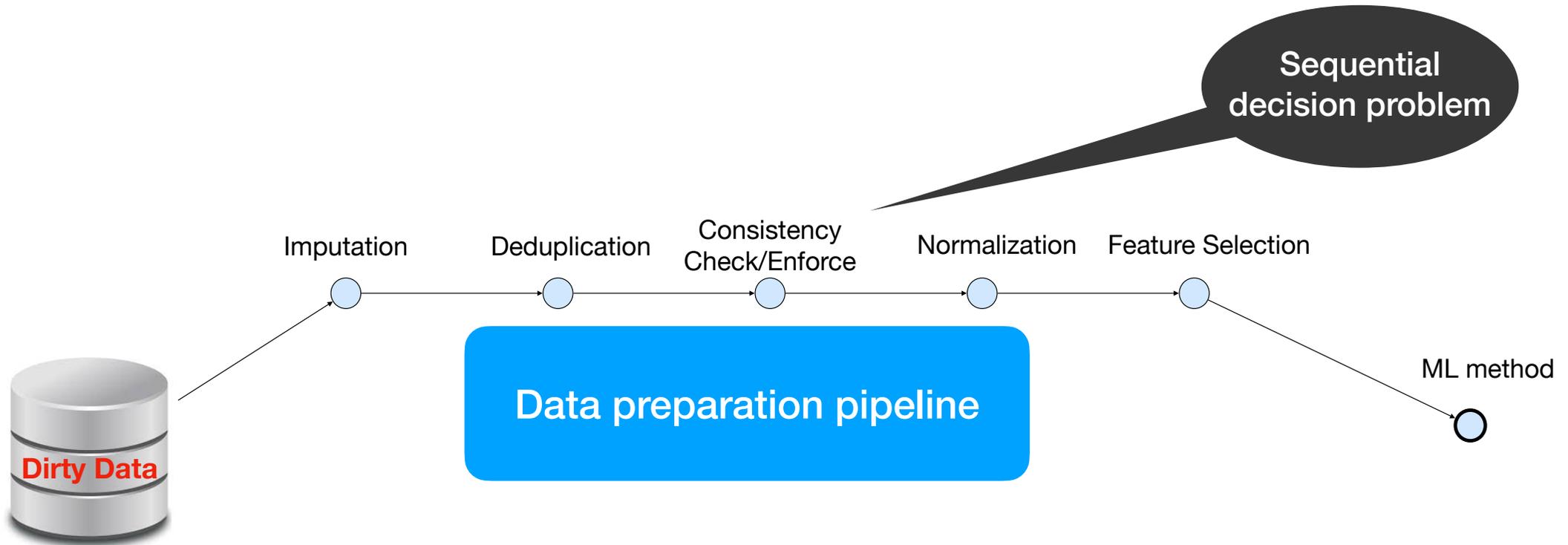
# Data preprocessing is challenging

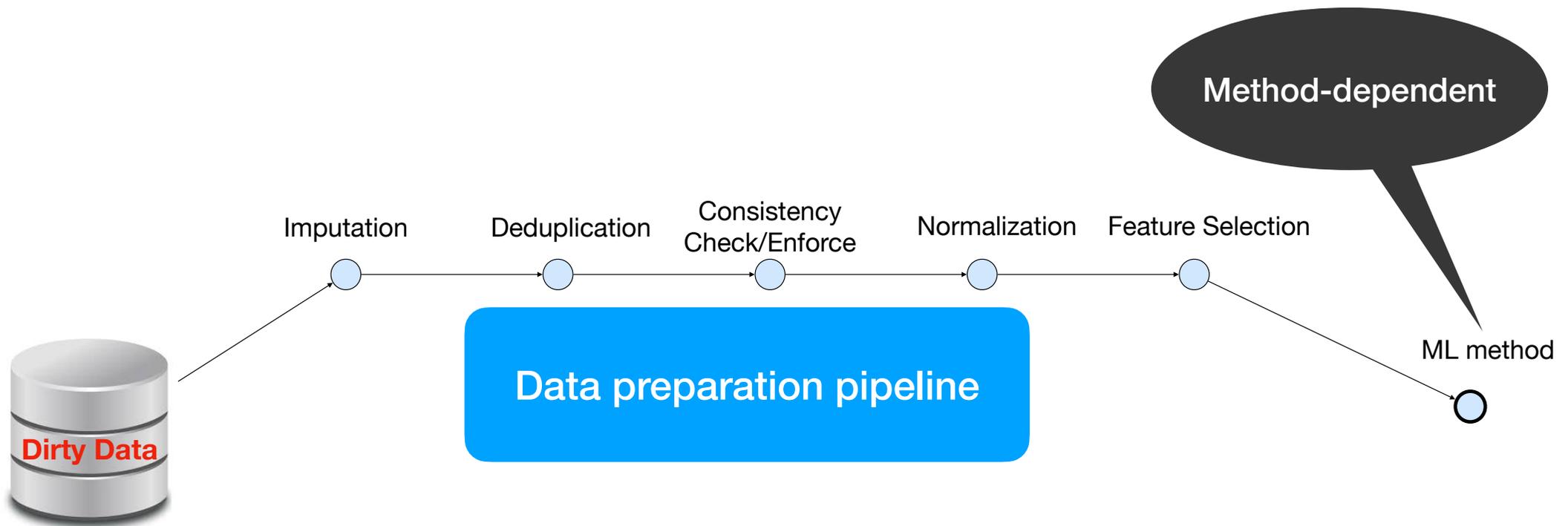**Dirty Data** → **Data preparation pipeline** → ML method

# Data preprocessing is challenging

Sequential decision problem

Imputation    Deduplication    Consistency Check/Enforce    Normalization    Feature Selection

Dirty Data

**Data preparation pipeline**

ML method

# Data preprocessing is challenging

**Method-dependent**

Imputation    Deduplication    Consistency Check/Enforce    Normalization    Feature Selection

**Data preparation pipeline**

ML method

Dirty Data

# Data preprocessing is challenging

**So many methods and parameter settings**

**Dirty Data**

**Imputation**
- Hot Deck
- MICE
- IRMI
- Median
- Mean
- Most Frequent
- K-NN
- …

**Deduplication**
- EditDistance
- Token-based
- N-grams
- FIFO
- Fusion
- …

**Consistency Check/Enforce**
- Rule-based
- FD-based
- Constraints
- Patterns
- …

**Normalization**
- Zscore—based
- Decimal scaling
- MinMax
- …

**Feature Selection**
- Missing ratio
- Linear correlation
- Model-based
- …

**ML method**

# Data preprocessing is challenging

Different orderings matter

Dirty Data

Imputation → Deduplication → Consistency Check/Enforce → Normalization → Feature Selection

Deduplication → Normalization → Imputation → Feature Selection

ML method

# Data preprocessing is challenging

# Data preprocessing is challenging

Selective processing of some parts of the dataset
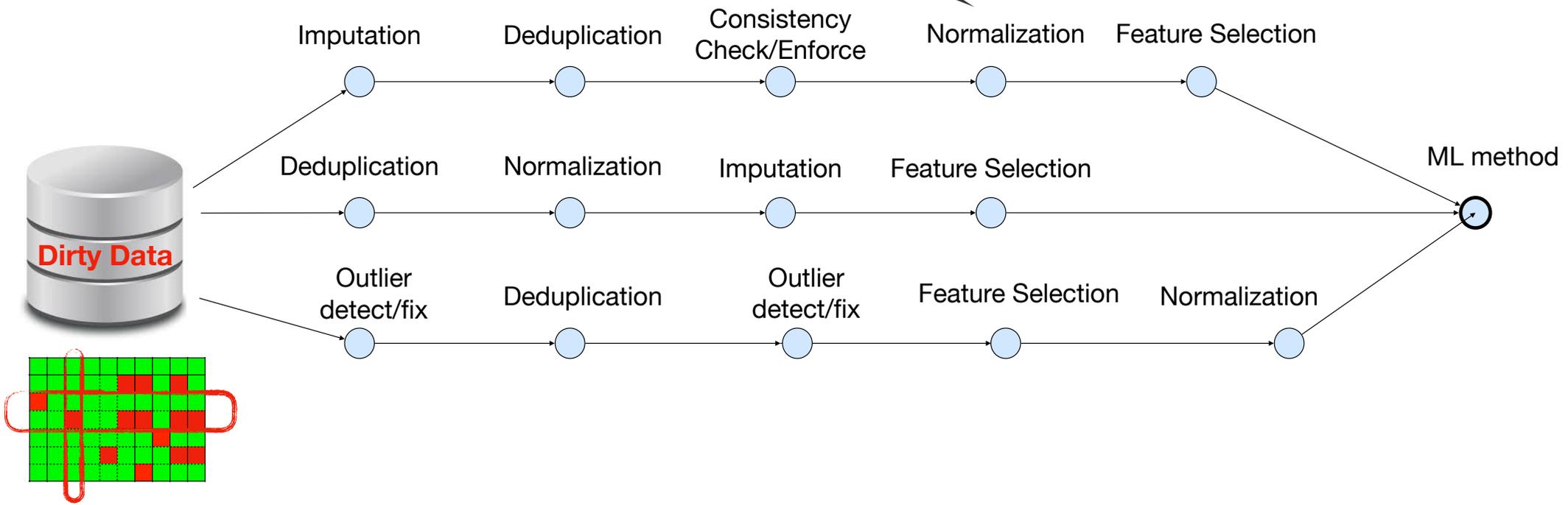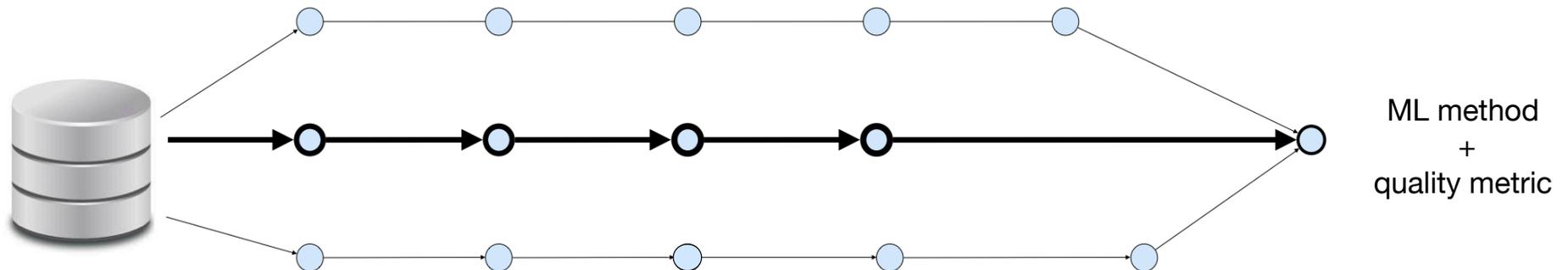
Dirty Data

| Imputation | Deduplication | Consistency Check/Enforce | Normalization | Feature Selection |

| Deduplication | Normalization | Imputation | Feature Selection |

ML method

| Outlier detect/fix | Deduplication | Outlier detect/fix | Feature Selection | Normalization |

# Data preprocessing is challenging



Dirty Data

Imputation → Deduplication → Consistency Check/Enforce → Normalization → Feature Selection

Deduplication → Normalization → Imputation → Feature Selection

Outlier detect/fix → Deduplication → Outlier detect/fix → Feature Selection → Normalization

ML method

**Patterns of glitches require specific data cleaning strategies**

[ICDE 2011]

43

# Data preprocessing is challenging

Infinite space of possible strategies

Dirty Data

| Imputation | Deduplication | Consistency Check/Enforce | Normalization | Feature Selection |

| Deduplication | Normalization | Imputation | Feature Selection |

| Outlier detect/fix | Deduplication | Outlier detect/fix | Feature Selection | Normalization |

ML method

# Optimization Problem

Can we help the user in composing the data preparation pipeline that maximizes the quality performance of the ML method ?



ML method
+
quality metric

# Optimization Problem

Can we help the user in composing the data preparation pipeline that maximizes the quality performance of the ML method ?

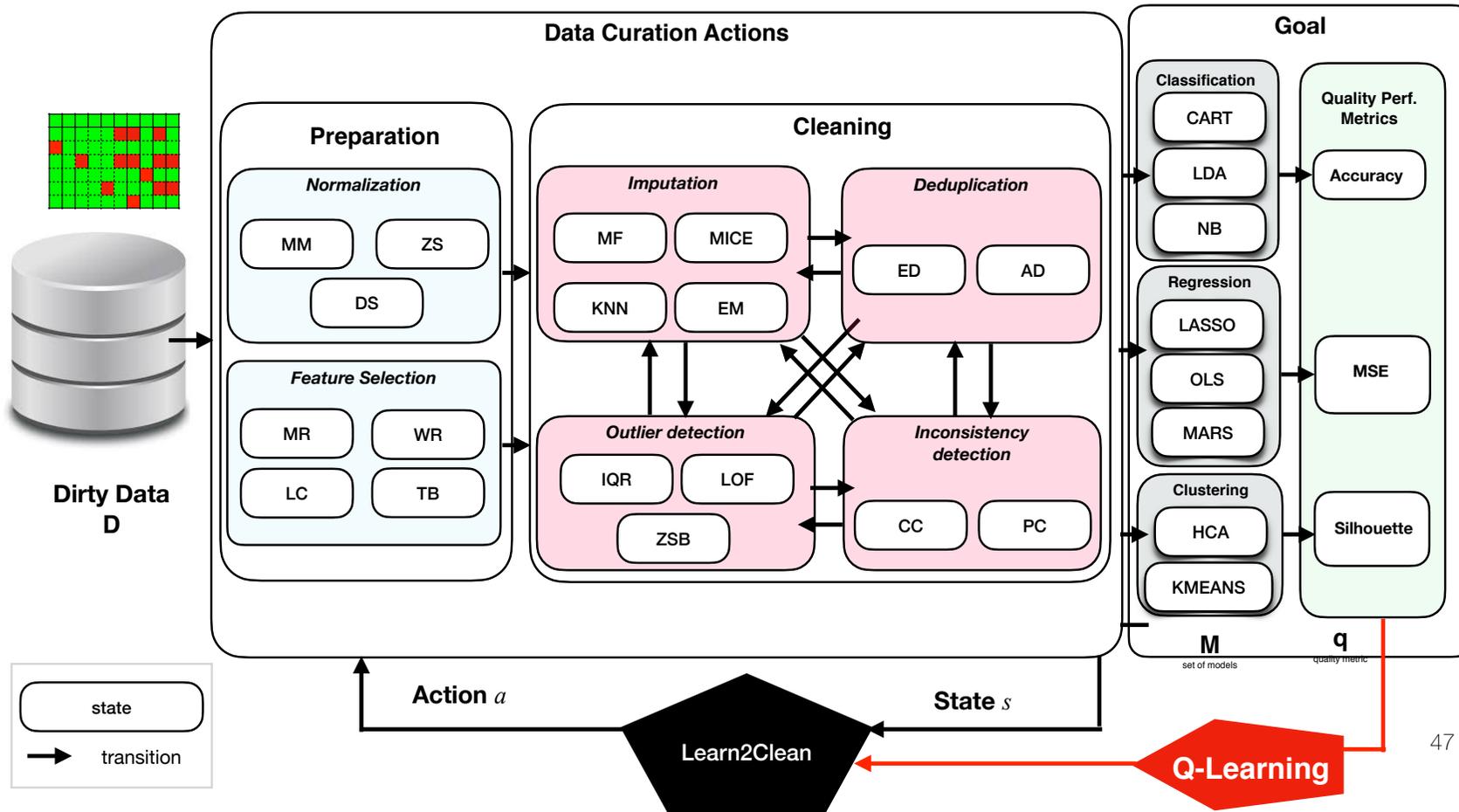No training example for "good" data cleaning

AutoML-like approach

Human-In-The Loop

Metric-dependent

No model a priori

ML method
+
quality metric

# First Solution: Learn2Clean

https://github.com/LaureBerti/Learn2Clean



**Data Curation Actions**

**Preparation**

*Normalization*
MM | ZS
DS

*Feature Selection*
MR | WR
LC | TB

**Cleaning**

*Imputation*
MF | MICE
KNN | EM

*Deduplication*
ED | AD

*Outlier detection*
IQR | LOF
ZSB

*Inconsistency detection*
CC | PC

**Dirty Data D**

**Goal**

**Classification**
CART
LDA
NB

**Regression**
LASSO
OLS
MARS

**Clustering**
HCA
KMEANS

**Quality Perf. Metrics**
Accuracy
MSE
Silhouette

**M** set of models
**q** quality metric

*AutoML-like approach for data Curation*

**Action** $a$

**State** $s$
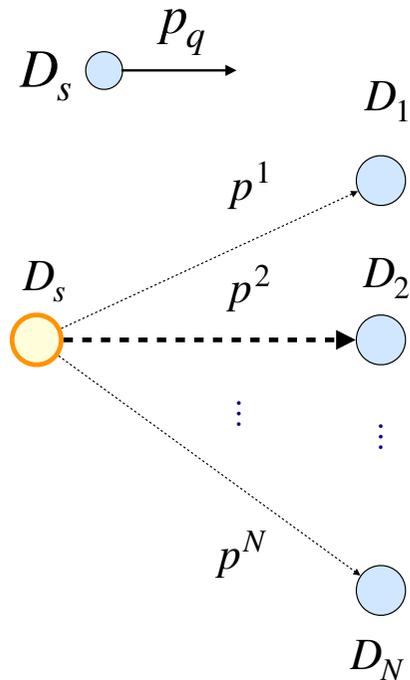
**Learn2Clean**

**Q-Learning**

state
→ transition

47

47

# Reinforcement Learning Framework

**Markov Decision Process**    State    Action    Transition    Reward    **Learn2Clean**

$$D_s \xrightarrow{p_q}$$

# Reinforcement Learning Framework

**Markov Decision Process**     State     Action     Transition     Reward        **Learn2Clean**

$$D_s \bigcirc \xrightarrow{p_q}$$

$$D_s$$
$$\bigcirc$$

# Reinforcement Learning Framework

**Markov Decision Process**  State  **Action**  Transition  Reward  **Learn2Clean**



$$D_s \xrightarrow{p_q}$$

$D_1$

$A$

$p^1$

$D_s$  $p^2$  $D_2$

$\vdots$  $\vdots$

$p^N$

$D_N$

# Reinforcement Learning Framework

**Markov Decision Process**  State  **Action**  Transition  Reward  **Learn2Clean**

$D_s$  $\xrightarrow{p_q}$

$A$

$D_s$ → $D_1$

$p^1$

$D_s$  $p^2$  $D_2$

$p^N$

$D_N$

| | |
|---|---|
| MICE<br>EM<br>KNN<br>MF | imputation |
| DS<br>MM<br>ZS | normalization |
| MR<br>WR<br>LC<br>TB | feature selection |
| ZSB<br>LOF<br>IQR | outlier detect/fix |
| CC<br>PC | consistency check/fix |
| AD<br>ED | duplicate detect/fix |
| LASSO or OLS or MARS for regression<br>HCA or KMEANS for clustering<br>CART or LDA or NB for classification | |

# Reinforcement Learning Framework

**Markov Decision Process**    State    Action    **Transition**    Reward    Learn2Clean
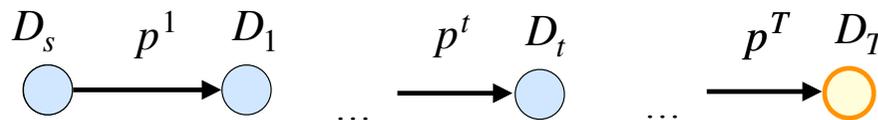
$$D_s \quad \xrightarrow{\quad p_q \quad}$$

$$D_s \quad \xrightarrow{\quad p^1 \quad} \quad D_1$$

# Reinforcement Learning Framework

**Markov Decision Process** $\boxed{\text{State}}$ $\boxed{\text{Action}}$ $\boxed{\textbf{Transition}}$ $\boxed{\text{Reward}}$ **Learn2Clean**

$$D_s \bigcirc \xrightarrow{p_q}$$

$$D_s \bigcirc \xrightarrow{p^1} D_1 \bigcirc \xrightarrow{p^2} D_2 \bigcirc$$

53

# Reinforcement Learning Framework

**Markov Decision Process**   State   Action   **Transition**   Reward   **Learn2Clean**

$$D_s \quad \circ \xrightarrow{\;p_q\;}$$

**Final State**

$$D_s \xrightarrow{\;p^1\;} D_1 \quad \dots \quad \xrightarrow{\;p^t\;} D_t \quad \dots \quad \xrightarrow{\;p^T\;} D_T$$

LASSO or OLS or MARS for regression
HCA or KMEANS for clustering
CART or LDA or NB for classification
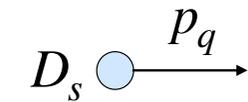
# Reinforcement Learning Framework

**Markov Decision Process** State Action Transition **Reward** **Learn2Clean**

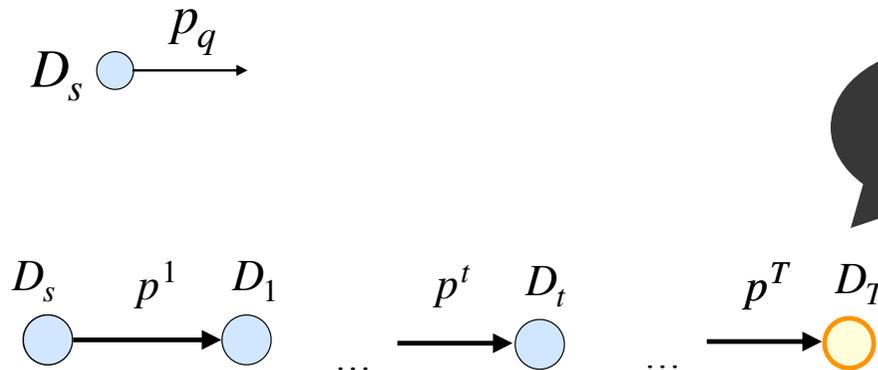$$D_s \xrightarrow{p_q}$$

**Final State**

*deterministic*

$D_s \xrightarrow{p^1} D_1 \quad \xrightarrow{p^t} D_t \quad \xrightarrow{p^T} D_T$

...    ...

MICE EM KNN MF DS MM ZS MR WR LC TB ZSB LOF IQR CC PC AD ED LASSO

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 100 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix}$$

55

# Reinforcement Learning Framework

**Markov Decision Process**     State     Action     Transition     Reward     **Learn2Clean**



$D_s$ $\xrightarrow{p_q}$

*deterministic*

**Final State**

$D_s$ $\xrightarrow{p^1}$ $D_1$ ... $\xrightarrow{p^t}$ $D_t$ ... $\xrightarrow{p^T}$ $D_T$     $R =$

MICE EM KNN MF DS MM ZS MR WR LC TB ZSB LOF IQR CC PC AD ED LASSO

LASSO or OLS or MARS for regression → MSE
HCA or KMEANS for clustering → Silhouette
CART or LDA or NB for classification → Accuracy

**Quality metric**

56

# Reinforcement Learning Framework

**Markov Decision Process**  State  Action  Transition  **Reward**  **Learn2Clean**

$$D_s \xrightarrow{p_q}$$

*deterministic*

**Final State**

MICE EM KNN MF DS MM ZS MR WR LC TB ZSB LOF IQR CC PC AD ED LASSO

$$D_s \xrightarrow{p^1} D_1 \quad \cdots \quad \xrightarrow{p^t} D_t \quad \cdots \quad \xrightarrow{p^T} D_T$$

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 100 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 \end{bmatrix}$$

LASSO or OLS or MARS for regression ⟶ MSE
HCA or KMEANS for clustering ⟶ Silhouette
CART or LDA or NB for classification ⟶ Accuracy
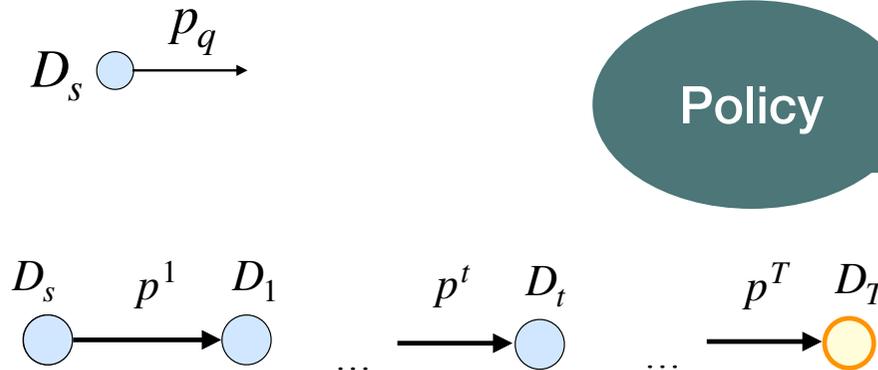
**Quality metric**

**Learn2Clean selects the sequence of preprocessing actions that maximizes the quality metric (or minimizes the error)**

# Reinforcement Learning Framework

**Markov Decision Process**   ( State )  ( Action )  ( Transition )  **Reward**   **Learn2Clean**

$$D_s \xrightarrow{p_q}$$

**Policy**

**Softmax action selection**

$$\pi = P(a \mid s) = \frac{e^{Q(s,a)/k}}{\sum_j e^{Q(s,a_j)/k}}$$

$$D_s \xrightarrow{p^1} D_1 \quad \dots \quad \xrightarrow{p^t} D_t \quad \dots \quad \xrightarrow{p^T} D_T$$

**Q-table**

*value iteration update*

$$Q^\pi(s,a) \leftarrow (1-\alpha) . Q(s,a) + \alpha . \left( \overbrace{R(s,a) + \gamma . \max_{a'} Q(s',a')}^{\text{learned value}} \right)$$

new value    learning rate    old value    reward    discount factor    optimal future value

# Experiments

## Datasets

| Name | # Att. | # Rows | Clustering | Regression | Classification |
|---|---|---|---|---|---|
| House Prices | 81 | 1.46k | ✓ | ✓ | ✓ |
| Google Playstore Users | 5 | 64.3k | ✓ | | |
| Google Playstore Apps | 13 | 10.8k | ✓ | | ✓ |

**Evaluation :** Silhouette for Clustering

MSE for Regression

Accuracy for Classification
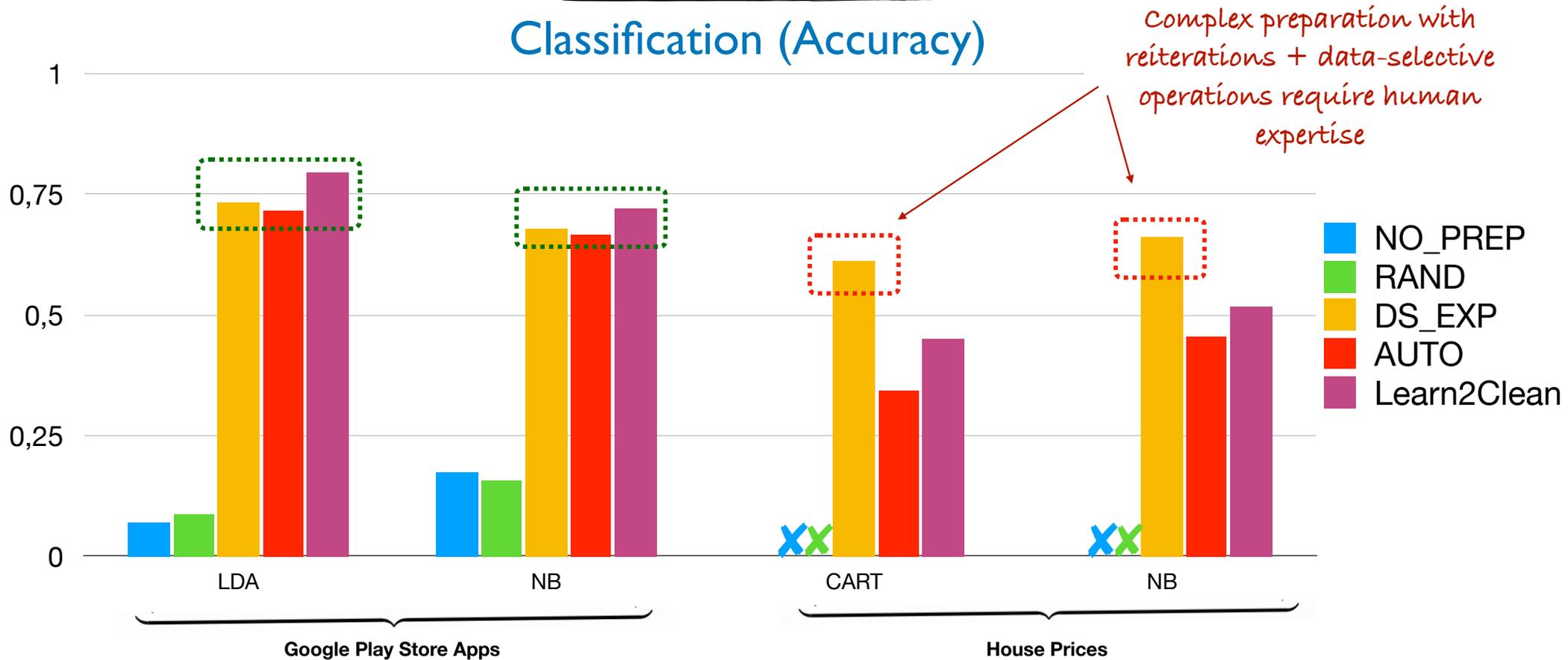
59

# Experimental Results

## Regression (MSE)



Legend:
- NO_PREP
- RAND
- DS_EXP
- AUTO
- Learn2Clean

House Prices

# Experimental Results

## Clustering (Silhouette)

# Experimental Results

## Classification (Accuracy)

Complex preparation with reiterations + data-selective operations require human expertise



Legend:
- NO_PREP
- RAND
- DS_EXP
- AUTO
- Learn2Clean

Google Play Store Apps: LDA, NB

House Prices: CART, NB

# HIL with Active Reward Learning

# Active Reward Learning

Goal: *learn from user feedbacks to adapt the rewards*

More likely to be chosen by the user

$\pi(a_t^i \mid s_t)$

$D_t^1$

0.05

$D_s$ $p^1$ $D_1$ ... $p^t$ $D_t$

$p_1^t$

$D_t^2$

$p_2^t$ 0.7

$p_{N_t}^t$

0.02

$D_t^{N_t}$

Learn2Clean
+
HIL

# Active Reward Learning

Goal: *learn from user feedbacks to adapt the rewards*

More likely to be chosen by the user

$\pi(a_t^i \mid s_t)$

$D_t^1$

$p_1^t$

$D_t^2$

$D_s$  $p^1$  $D_1$  $p^t$  $D_t$  $p_2^t$  ✖

... 

$p_{N_t}^t$  ✖

$D_t^{N_t}$

0.05

0.7

0.02

65

# Active Reward Learning

Learn2Clean
+
HIL

Goal: **learn from user feedbacks to adapt the rewards**

More likely to be chosen by the user

$\tilde{\pi}(a_t^i \mid s_t)$    $\pi(a_t^i \mid s_t)$

$D_t^1$

$p_1^t$

0.9    0.05

$D_t^2$

$D_s$    $p^1$    $D_1$    $p^t$    $D_t$    $p_2^t$

0    Force    0.7

exploration

$p_{N_t}^t$

0    0.02

$D_t^{N_t}$

# Ongoing work

- New version of Learn2Clean with deep RL agents
- Combine AutoML, AutoCuration, and HIL
- Learn better reward functions
- Extend the library of ML and data preparation methods
- Extend experiments with more intricate data glitches and various glitch distributions





Code: https://github.com/LaureBerti/Learn2Clean

# Concluding Remarks

- ML crucially needs principled data curation and preparation, adequate tooling, and user assistance
- The impact of data preprocessing variability is largely underestimated in ML
- Many data preprocessing tasks require seamless integration of <u>Human-in-the-Loop</u> and <u>automated ML-based</u> solutions
- Perfect timing for many R&D opportunities:
  - Manage and orchestrate human/machine resources
  - Challenge and transfer research ideas to operational and very large-scale contexts

Thank you!