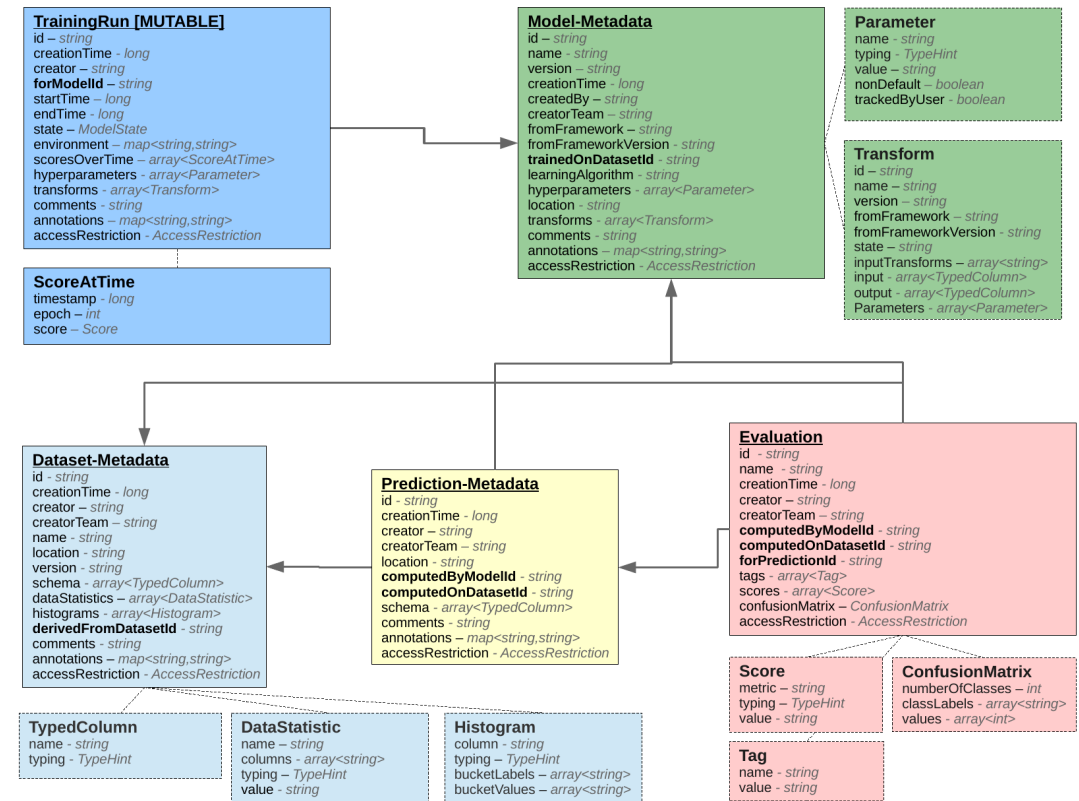# Optimization of Machine Learning Workloads with Experiment Databases

Behrouz Derakhshan

Behrouz.derakhshan@dfki.de

# Experiment Database (ED)

- Contains meta-data of machine learning experiments about the datasets, models, transformations, and evaluations

- Based on [1]

- Workload Optimizations using ED:
  - Reuse of existing models
  - Materialization of transformed datasets
  - Multi-query optimization
  - Warm-starting
  - Efficient data transformation
  - Efficient hyper-parameter tuning
  - Efficient multi-model training

**TrainingRun [MUTABLE]**
id – *string*
creationTime - *long*
creator – *string*
**forModelId** – *string*
startTime - *long*
endTime - *long*
state – *ModelState*
environment – *map<string,string>*
scoresOverTime – *array<ScoreAtTime>*
hyperparameters - *array<Parameter>*
transforms - *array<Transform>*
comments - *string*
annotations – *map<string,string>*
accessRestriction - *AccessRestriction*

**ScoreAtTime**
timestamp - *long*
epoch – *int*
score – *Score*

**Model-Metadata**
id – *string*
name - *string*
version - *string*
creationTime - *long*
createdBy – *string*
creatorTeam – *string*
fromFramework – *string*
fromFrameworkVersion - *string*
**trainedOnDatasetId** - *string*
learningAlgorithm - *string*
hyperparameters - *array<Parameter>*
location - *string*
transforms - *array<Transform>*
comments - *string*
annotations – *map<string,string>*
accessRestriction - *AccessRestriction*

**Parameter**
name - *string*
typing - *TypeHint*
value – *string*
nonDefault – *boolean*
trackedByUser - *boolean*

**Transform**
id – *string*
name – *string*
version – *string*
fromFramework – *string*
fromFrameworkVersion - *string*
state – *string*
inputTransforms – *array<string>*
input - *array<TypedColumn>*
output - *array<TypedColumn>*
Parameters - *array<Parameter>*

**Dataset-Metadata**
id - *string*
creationTime - *long*
creator – *string*
creatorTeam – *string*
name - *string*
location - *string*
version - *string*
schema - *array<TypedColumn>*
dataStatistics – *array<DataStatistic>*
histograms - *array<Histogram>*
**derivedFromDatasetId** - *string*
comments - *string*
annotations – *map<string,string>*
accessRestriction - *AccessRestriction*

**Prediction-Metadata**
id - *string*
creationTime - *long*
creator – *string*
creatorTeam – *string*
location - *string*
**computedByModelId** - *string*
**computedOnDatasetId** - *string*
schema - *array<TypedColumn>*
comments - *string*
annotations – *map<string,string>*
accessRestriction - *AccessRestriction*

**Evaluation**
id - *string*
name - *string*
creationTime - *long*
creator – *string*
creatorTeam – *string*
**computedByModelId** - *string*
**computedOnDatasetId** - *string*
**forPredictionId** - *string*
tags - *array<Tag>*
scores - *array<Score>*
confusionMatrix – *ConfusionMatrix*
accessRestriction - *AccessRestriction*

**Score**
metric – *string*
typing – *TypeHint*
value - *string*

**ConfusionMatrix**
numberOfClasses – *int*
classLabels - *array<string>*
values - *array<int>*

**Tag**
name - *string*
value - *string*

**TypedColumn**
name - *string*
typing - *TypeHint*

**DataStatistic**
name – *string*
columns - *array<string>*
typing – *TypeHint*
value - *string*

**Histogram**
column – *string*
typing – *TypeHint*
bucketLabels – *array<string>*
bucketValues – *array<string>*

Database schema for storing declarative description of machine learning experiments [1]

# Categories of Machine Learning Workloads

## 1. Interactive batch workload

- Description:
  - Multi-user (from tens to thousands of users)
  - Primarily consists of repeated data preprocessing and data transformations
  - Aggressive Hyper-parameter tuning and model evaluation

- Optimizations using ED:
  - Reuse
  - Materialization
  - Multi-query optimization
  - Warm-starting
  - Decreasing the search space of hyper-parameters

## 2. Incremental Training Workload

- Description:
  - Multi-user (tens of users)
  - Incremental improvement of existing data processing pipelines and ML models
  - Retrain models and pipelines (typically daily)

- Optimizations using ED:
  - Speed up in data transformation
  - Decreasing the search space of hyper-parameters
  - Multi-model training
  - Warm-starting

## 3. Continuous Training Workload

- Description:
  - Real-time (or near real-time) data processing
  - Online ML models

- Optimizations using ED:
  - Speed up in data preprocessing
  - Multi-model training
  - Fast detection of model quality degradation

# Plan for next Paper

- VLDB 2018 (Deadline: 1$^{st}$ of March)

- Focus on Workload 1 (Interactive batch)

- Plan:
  - Develop a simple database based on [1]
  - Experiments:
    - Focus only on a few frameworks (scikit-learn, R mlr package)
    - Popular Kaggle competitions [3], OpenML datasets and tasks [2]
    - I expect a reduction in the processing time and the development time

- I currently have a paper under review for SIGMOD2018. Depending on the result of the review, this work may have to be pushed back to SIGMOD2019/VLDB2018

- Early Results:
  - The Most popular pipeline in OpenML (scikit-learn) consists of:
    - Missing Value Imputer
    - Dimensionality reduction using PCA
    - Random Forest Classifier
  - It is repeatedly executed on 100 different tasks (average of 9 times on each task)
  - Figure below shows that a simple reuse can save about 2 hours of processing time



Total vs Single execution time

This figure may include executions performed by bots, therefore, we cannot reach a conclusion about the reduction on the development time

# References

[1] *Sebastian Schelter, Joos-Hendrik Böse, Johannes Kirschnick, Thoralf Klein, Stephan Seufert*, **Automatically Tracking Metadata and Provenance of Machine Learning Experiments**, Machine Learning Systems workshop at the conference on Neural Information Processing Systems (NIPS) 2017

[2] https://www.openml.org/

[3] https://www.kaggle.com/