

AUTOMATED DOCUMENTATION OF END-TO-END EXPERIMENTS IN DATA SCIENCE

Sergey Redyuk, TU Berlin, Germany sergey.redyuk@tu-berlin.de
 supervised by Volker Markl (TU Berlin) and Sebastian Schelter (NYU)

Abstract

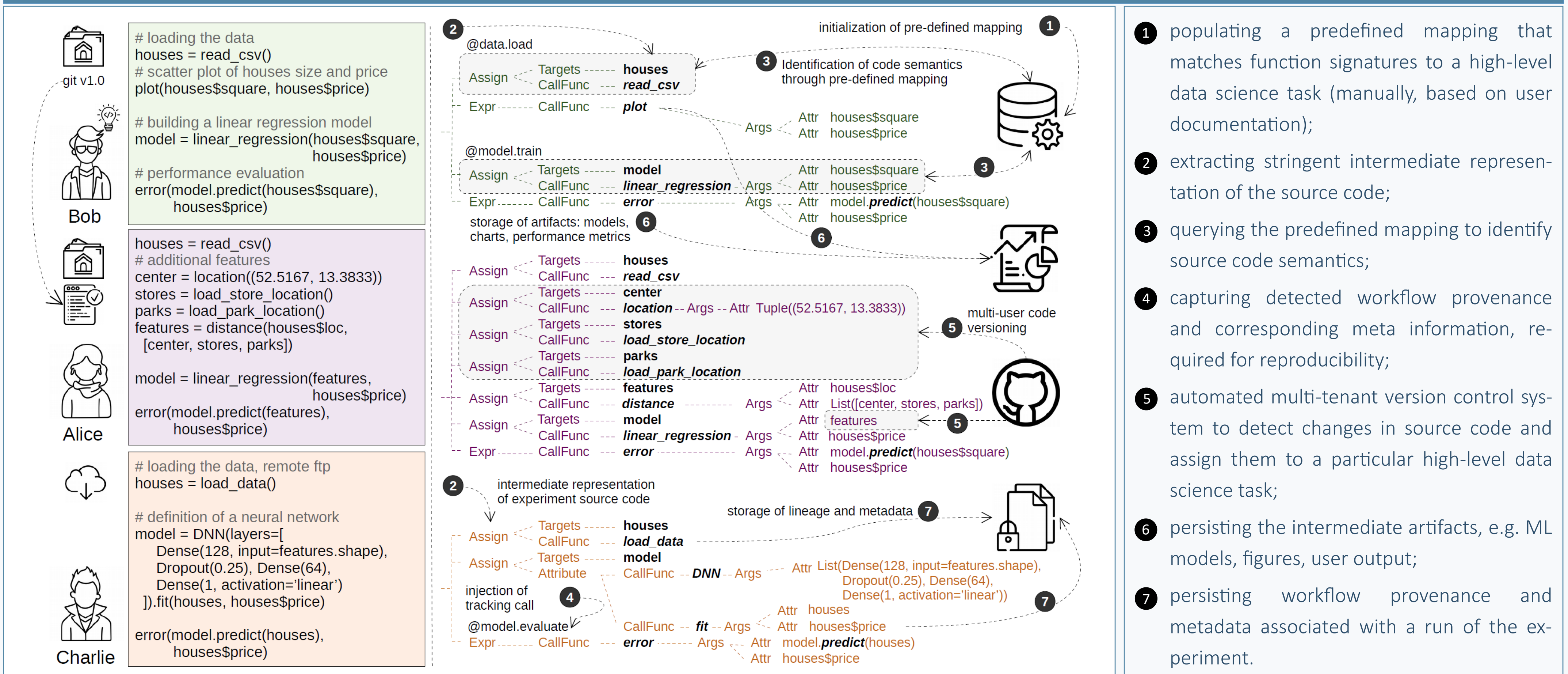
Data-oriented experiments are **hard to reproduce** due to rapid undocumented changes, abundance and inconsistency of tools and frameworks, multi-tenant environment

Goal: end-to-end tracking of workflow provenance and meta information, to achieve reproducibility, sharing and reuse of intermediate artifacts across teams and domains

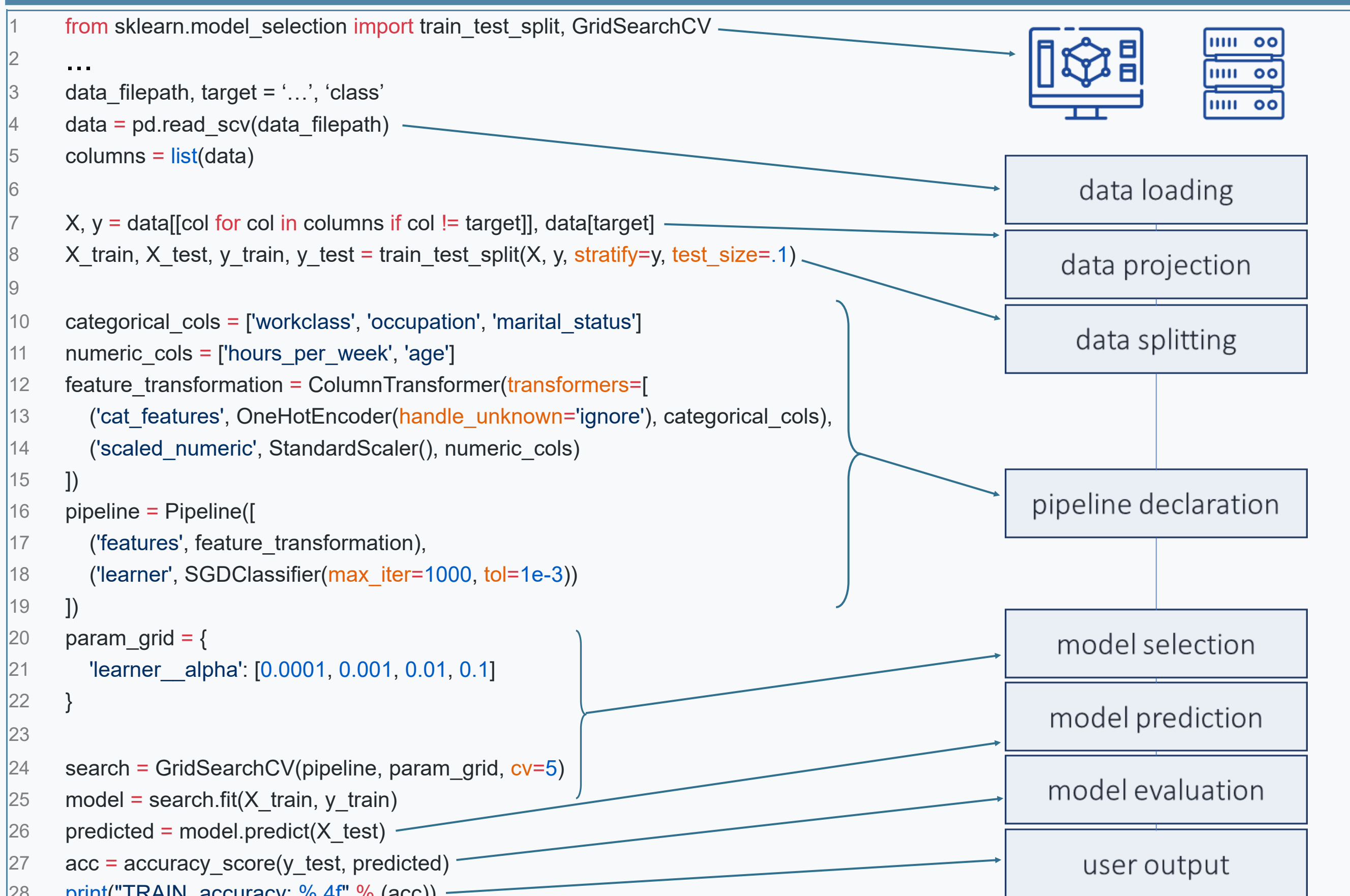
Contributions

- automated tracking of workflow provenance and metadata at runtime via source code analysis
- extraction of a declarative high-level representation of a data science experiment
- design and population of an experiment database to enable search and reuse capabilities over previous experiments

Use Case



Open Questions and Challenges



Evaluation

- manually annotated example scripts where ground truth is available, e.g. logical structure of the experiment, ML model declaration, hyperparameters etc.
- example scripts collected from public repositories where ground truth can be partially inferred, e.g. jupyter notebooks that are logically splitted into blocks
- metrics: achieved reproducibility (binary), overhead as a diff in execution time

Semantics

- difficult to infer, e.g. low-level implementations
- the task on arbitrary source code is so far infeasible
- experiment source code is a combination of high-level APIs (e.g. pandas ecosystem, keras) and low-level imperative statements
- solutions: (i) frequent pattern mining, and (ii) evaluation of control flow



We would like to acknowledge the support of the Helmholtz Einstein International Berlin Research School in Data Science (HEIBRIDS), Max Delbrück Center for Molecular Medicine (MDC), TU Berlin, and the Berlin Center for Machine Learning (BZML) 01S18037A. We thank Sebastian Schelter (NYU) and Volker Markl (TU Berlin) for guidance and valuable contributions.

