

mlinspect: Inspect ML Pipelines in Python in the form of a DAG

Stefan Grafberger

Technical University Munich

November 2, 2020

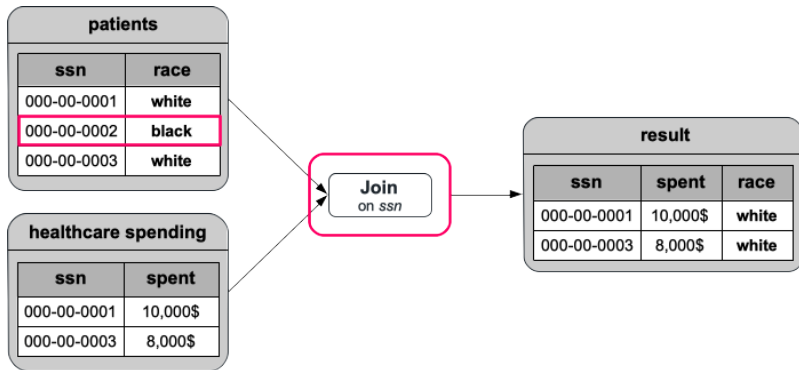
Motivation

- ML is everywhere and is increasingly used to automate decisions that impact peoples' lives, some examples:
 - loan applications
 - medical diagnosis
 - job applications
- ML Systems can be very brittle and there are many problems with the fairness, transparency and accountability of these systems
- Issues can be hard to detect without manual and time-intensive scrutiny

Issue Sources

- Data
 - ⇒ Data quality issues, as well as data over- and under-representation, can heavily impact the results
- ML Model
 - ⇒ Lots of research from the ML community on this, but mostly focuses on the adaption of learning algorithms on static datasets
- Preprocessing code
 - Input data for ML applications has to be integrated, preprocessed and cleaned first
 - Creating a good ML preprocessing pipeline is difficult and requires a lot of ML expertise
 - ⇒ We need to look at the pipeline as a whole, not just at data and model!

Example



Operations like Selections, Joins and Missing Value imputation can introduce data distribution issues

⇒ Similar to code inspections in modern IDEs

- Can automatically detect issues like statistical bias and provides linting for best-practices
- Can also help with debugging
- Works on existing pipeline code using popular libraries without requiring modifications
- Negligible performance overhead

Embedding vectors
may not be available
for rare names!

Declarative inspection of preprocessing pipeline

```
PipelineInspector
.on_pipeline('health.py')
.no_bias_introduced_for(
    ['age_group', 'race'])
.no_illegal_features()
.no_missing_embeddings()
.verify()
```

6 / 6