

Big Data and Cloud Computing, 18-19

Inês Dutra
DCC-FCUP

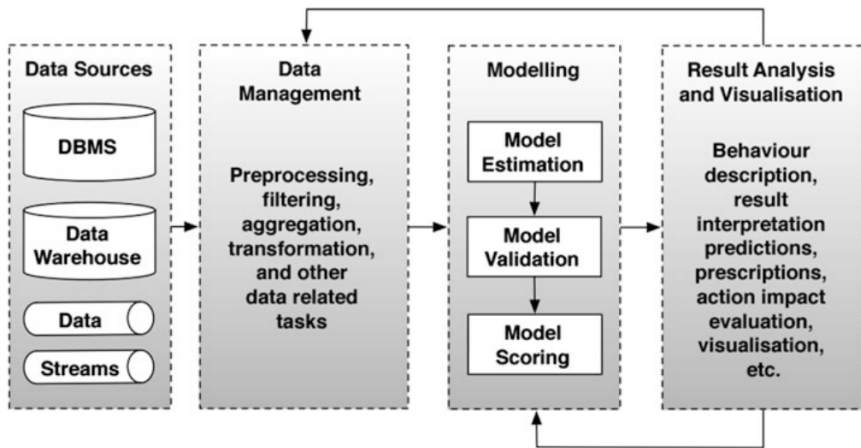
ines@dcc.fc.up.pt (gabinete: 1.31)

March 26, 2019

Data Mining and Machine Learning: recap

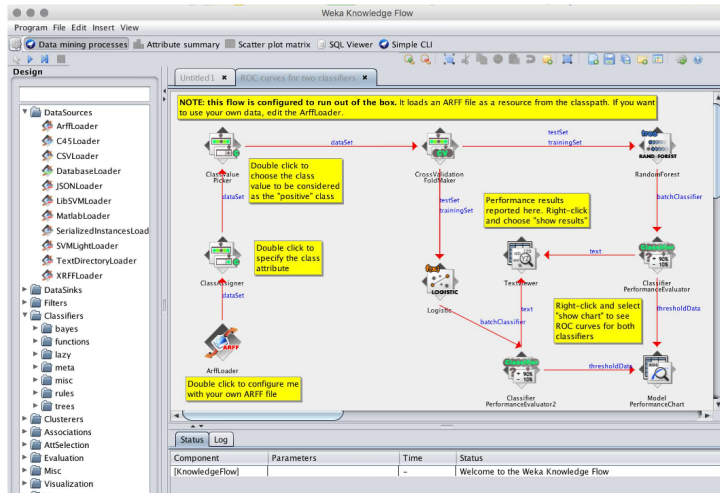
- Workflow (Dataflow - Knowledgeflow):
 - ▶ Data preprocessing
 - transformation: normalization, standardization, averaging, median, denoising, filtering
 - preparation: depends on the task, algorithm, package or library being used
 - ▶ Machine learning task, algorithm
 - ▶ Validation: cross-validation, bootstrapping
- Tools: WEKA, RapidMiner, Taverna, Condor DAGMan, Pegasus, Google Dataflow, Google Composer (Apache Airflow)

Data Mining and Machine Learning: workflow



Assunção, M.D., Calheiros, R.N., Bianchi, S., Netto, M.A.S., Buyya, R.: Big data computing and clouds: trends and future directions. *J. Parallel Distrib. Comput.* 79–80, 3–15 (2015)

Example of workflow in WEKA

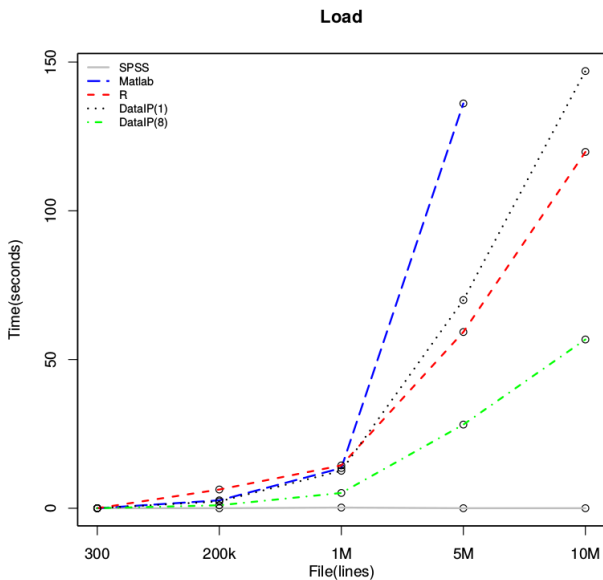


java -jar weka.jar ⇒ KnowledgeFlow

Scalability

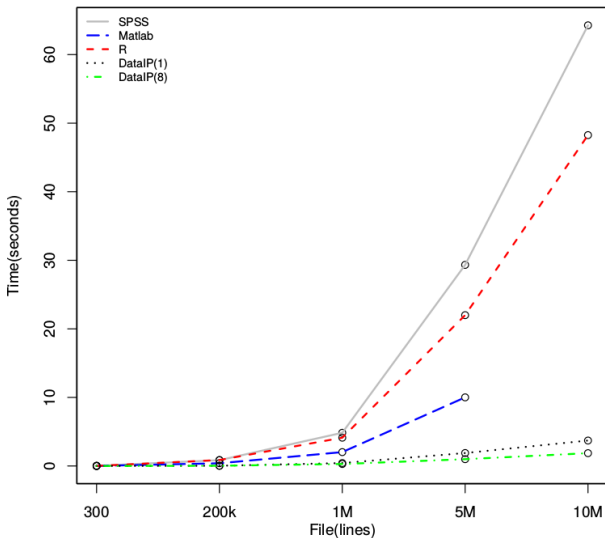
- Computational resources: memory, CPU, I/O, storage
- I/O:
 - ▶ **Experiment 1:** SPSS, MatLab, R and DataIP (in-house implementation)
 - dataset of patients, originally 200K entries, 6 numeric variables without nulls
 - varying sizes: 300, 200k, 1M, 5M, 10M
 - ▶ **Experiment 2:** job that needs to fetch data files from a remote site

Scalability: Experiment 1, I/O



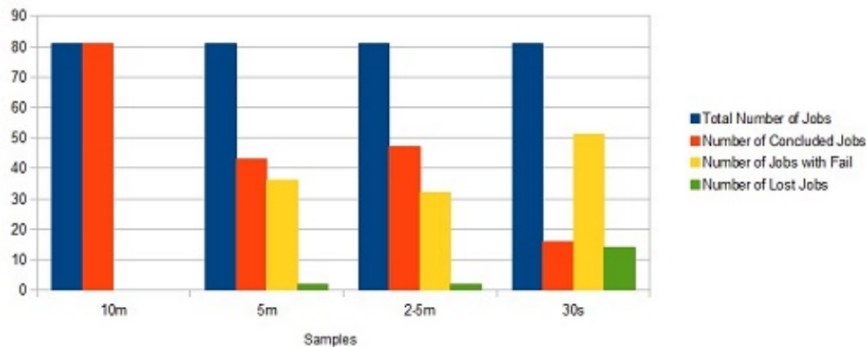
Scalability: Experiment 1, simple computing: summary

Summary



Scalability: Experiment 2, file transfer

Jobs x Samples



Scalability

- Alternatives
 - ▶ break file in multiple smaller files that can be read in parallel: useful if the reading can be done in parallel
 - ▶ undersampling: need to be careful about data distribution
 - ▶ use of specific hardware and software: distributed disks, distributed file systems, distributed databases, in-memory databases, parallel and distributed software
 - ▶ work with compressed files: zip, parquet, CSR, CSR5 (for sparse matrices) etc