

SCIENCE PASSION TECHNOLOGY

Data Integration and Analysis 01 Introduction and Overview

Matthias Boehm

Graz University of Technology, Austria Computer Science and Biomedical Engineering Institute of Interactive Systems and Data Science BMK endowed chair for Data Management









Announcements/Org

- #1 Video Recording
 - Link in TUbe & TeachCenter (lectures will be public)
 - Optional attendance (independent of COVID)
 - Hybrid, in-person but video-recorded lectures
 - HS i5 + Webex: <u>https://tugraz.webex.com/meet/m.boehm</u>
- #2 COVID-19 Precautions (HS i5)
 - Room capacity: 24/48 (green/yellow), 12/48 (orange/red)
 - TC lecture registrations (limited capacity, contact tracing)
- **#3 Course Registration** (as of Oct 07)
 - Data Integration and Large-Scale Analysis (DIA)

WS20/21: **96 (2)** WS21/22: **122 (4)**



TUbe



Announcements/Org, cont.

- #4 Study Abroad Fair
 - International Days 2021
 - Oct 19 21, 2021
 - Virtual presentations, drop-in café
 - https://tu4u.tugraz.at/studierende/ mein-auslandsaufenthalt/ informationsveranstaltungen/international-days-2021/

#5 Open Data-Science Positions

- REDWAVE (ML engineer for recycling)
- Know-Center (data/ML engineer for automotive)













Agenda

- Data Management Group
- Course Organization
- Course Motivation and Goals
- Course Outline and Projects
- Excursus: Apache SystemDS





Data Management Group

https://damslab.github.io/





About Me

- **09/2018 TU Graz**, Austria
 - BMK endowed chair for data management
 - Data management for data science

(ML systems internals, end-to-end data science lifecycle)





Center

- 2012-2018 IBM Research Almaden, USA
 - Declarative large-scale machine learning
 - Optimizer and runtime of Apache SystemML
- 2011 PhD TU Dresden, Germany
 - Cost-based optimization of integration flows
 - Systems support for time series forecasting
 - In-memory indexing and query processing











Data Management Courses







Course Organization



Basic Course Organization

- Staff
 - Lecturer: Univ.-Prof. Dr.-Ing. Matthias Boehm, ISDS
 - Assistants: M.Sc. Shafaq Siddiqi, M.Sc. Sebastian Baunsgaard

Language

- Lectures and slides: English
- Communication and examination: English/German

Course Format

- VU 2/1, 5 ECTS (2x 1.5 ECTS + 1x 2 ECTS), bachelor/master
- Weekly lectures (Fri 3pm, including Q&A), attendance optional
- Mandatory exercises or programming project (2 ECTS)
- Recommended papers for additional reading on your own

Prerequisites

- Preferred: course Data Management / Databases is very good start
- Sufficient: basic understanding of SQL / RA (or willingness to fill gaps)
- Basic programming skills (Python, R, Java, C++)





9

¹⁰ Course Logistics

- Website
 - https://mboehm7.github.io/teaching/ws2122_dia/index.htm
 - All course material (lecture slides) and dates
- Video Recording Lectures (TUbe)?



Communication

- Informal language (first name is fine)
- Please, immediate feedback (unclear content, missing background)
- Newsgroup: N/A email is fine, summarized in following lectures
- Office hours: by appointment or after lecture
- Exam
 - Completed exercises or project (checked by me/staff)
 - Final written exam (oral exam if <25 students take the exam)
 - Grading (30% project/exercises completion, 70% exam)



Course Logistics, cont.

Course Applicability

- Bachelor programs computer science (CS), as well as software engineering and management (SEM)
- Master programs computer science (CS), as well as software engineering and management (SEM)
 - Catalog Data Science: compulsory course in major/minor
- Free subject course in any other study program or university





Course Motivation and Goals



Data Sources and Heterogeneity

- Terminology
 - Integration (Latin integer = whole): consolidation of data objects / sources
 - Homogeneity (Greek homo/homoios = same): similarity
 - Heterogeneity: dissimilarity, different representation / meaning

Heterogeneous IT Infrastructure

- Common enterprise IT infrastructure contains >100s of heterogeneous and distributed systems and applications
- E.g., health care data management: 20 120 systems

Multi-Modal Data (example health care)

- Structured patient data, patient records incl. prescribed drugs
- Knowledge base drug APIs (active pharmaceutical ingredients) + interactions
- Doctor notes (text), diagnostic codes, outcomes
- Radiology images (e.g., MRI scans), patient videos
- Time series (e.g., EEG, ECoG, heart rate, blood pressure)

















The 80% Argument

- Data Sourcing Effort
 - Data scientists spend 80-90% time on finding, integrating, cleaning datasets

[Michael Stonebraker, Ihab F. Ilyas: Data Integration: The Current Status and the Way Forward. IEEE Data Eng. Bull. 41(2) (2018)]

-	-	12	2	1
局			817	5
		1075		10
				6
10				
-03				-

Technical Debts in ML Systems



- Glue code, pipeline jungles, dead code paths
- Plain-old-data types (arrays), multiple languages, prototypes
- Abstraction and configuration debts
- Data testing, reproducibility, process management, and cultural debts







Horizontal Integration (e.g., EAI)

706.520 Data Integration and Large-Scale Analysis – 01 Introduction and Overview Matthias Boehm, Graz University of Technology, WS 2021/22





Course Goals

- Common Data and System Characteristics
 - Heterogeneous data sources and formats, often distributed
 - Large data collections → distributed data storage and analysis
- #1 Major data integration architectures
- #2 Key techniques for data integration and cleaning
- #3 Methods for large-scale data storage and analysis





Course Outline and Projects



Part A: Data Integration and Preparation

Data Integration Architectures

- 01 Introduction and Overview [Oct 08]
- 02 Data Warehousing, ETL, and SQL/OLAP [Oct 15]
- O3 Message-oriented Middleware, EAI, and Replication [Oct 22]

Key Integration Techniques

- 04 Schema Matching and Mapping [Oct 29]
- 05 Entity Linking and Deduplication [Nov 05]
- 06 Data Cleaning and Data Fusion [Nov 12, Shafaq]
- 07 Data Provenance and Blockchain [Nov 19]







Part B: Large-Scale Data Management & Analysis

Cloud Computing

- 08 Cloud Computing Foundations [Nov 26]
- O9 Cloud Resource Management and Scheduling [Dec 03]
- 10 Distributed Data Storage [Dec 10]

Large-Scale Data Analysis

- 11 Distributed, Data-Parallel Computation [Jan 07]
- **12 Distributed Stream Processing** [Jan 14]
- 13 Distributed Machine Learning Systems [Jan 21]
- 14 Q&A and exam preparation [Jan 21]
 - https://mboehm7.github.io/teaching/ws2021_dia/ExamDIA_v1.pdf
 - https://mboehm7.github.io/teaching/ws2021_dia/ExamDIA_v2.pdf





Overview Projects or Exercises

- Team
 - 1-3 person teams (w/ clearly separated responsibilities)
 - In exceptions also larger teams (e.g., Data Cleaning Benchmark)
- Objectives
 - Non-trivial programming project in DIA context (2 ECTS → 50 hours)
 - Preferred: Open source contribution to Apache SystemDS <u>https://github.com/apache/systemds</u> (from HW to high-level scripting)
 - Alternative Exercise: Data engineering and ML pipeline
 - Data integration and preparation of multi-modal data sources
 - Data cleaning and ML model training and evaluation
- Timeline
 - Oct 22: List of projects proposals / exercise description
 - Oct 29: Binding project/exercise selection
 - Jan 21: Final project/exercise deadline





Apache SystemDS: A Declarative ML System for the End-to-End Data Science Lifecycle

Background and System Architecture https://github.com/apache/systemds





Data Science Lifecycle

23



What is an ML System?



Data Science Lifecycle

24

Landscape of ML Systems

- Existing ML Systems
 - #1 Numerical computing frameworks
 - #2 ML Algorithm libraries (local, large-scale)
 - #3 Linear algebra ML systems (large-scale)
 - #4 Deep neural network (DNN) frameworks
 - #5 Model management, and deployment
- Exploratory Data-Science Lifecycle
 - Open-ended problems w/ underspecified objectives
 - Hypotheses, data integration, run analytics
 - Unknown value → lack of system infrastructure
 → Redundancy of manual efforts and computation
- Data Preparation Problem
 - **80% Argument:** 80-90% time for finding, integrating, cleaning data
 - Diversity of tools → boundary crossing, lack of optimization



"Take these datasets and show value or competitive advantage"

[DEBull 201	L8]	
data	A starts barren artik hard Martin Artikar Artikar artikar Artik	

[NIPS 2015]



Data-centric View:









Example: Linear Regression Conjugate Gradient

Note: #1 Data Independence #2 Implementation- Agnostic Operations	1: 2: 3: 4:	<pre>X = read(\$1); # n x m matrix y = read(\$2); # n x 1 vector maxi = 50; lambda = 0.001; intercept = \$3;</pre>	Read matrices from HDFS/S3
0	5: 6: 7:	<pre> r = -(t(X) %*% y); norm_r2 = sum(r * r); p = -r;</pre>	Compute initial gradient
Compute conjugate gradient	8: 9: 10: 11: 12:	<pre>w = matrix(0, ncol(X), 1); i = 0; while(i<maxi &="" norm_r2="">norm_r2_trgt) { q = (t(X) %*% (X %*% p))+lambda*p alpha = norm_r2 / sum(p * q);</maxi></pre>	; Compute
Update model and residuals	13: 14: 15: 16: 17: 18:	<pre>w = w + alpha * p; old_norm_r2 = norm_r2; r = r + alpha * q; norm_r2 = sum(r * r); beta = norm_r2 / old_norm_r2; p = -r + beta * p; i = i + 1;</pre>	step size
	19: 20:	<pre>} write(w, \$4, format="text");</pre>	of Concerns"



Apache SystemML/SystemDS





Cluster Config:

Basic HOP and LOP DAG Compilation

LinregDS (Direct Solve)





➔ Hybrid Runtime Plans:

- Size propagation / memory estimates
- Integrated CP / Spark runtime
- Dynamic recompilation during runtime

Distributed Matrices

- Fixed-size (squared) matrix blocks
- Data-parallel operations

Apache SystemDS Design

- Objectives
 - Effective and efficient data preparation, ML, and model debugging at scale
 - High-level abstractions for different lifecycle tasks and users
- #1 Based on DSL for ML Training/Scoring
 - Hierarchy of abstractions for DS tasks
 - ML-based SotA, interleaved, performance



Apache SystemML (since 2010)

→ Apache SystemDS (07/2020)

→ SystemDS (09/2018)

- #2 Hybrid Runtime Plans and Optimizing Compiler
 - System infrastructure for diversity of algorithm classes
 - Different parallelization strategies and new architectures (Federated ML)
 - Abstractions → redundancy → automatic optimization
- #3 Data Model: Heterogeneous Tensors
 - Data integration/prep requires generic data model



706.520 Data Integration and Large-Scale Analysis – 01 Introduction and Overview Matthias Boehm, Graz University of Technology, WS 2021/22



Language Abstractions and APIs, cont.

Example: Stepwise Linear Regression











[M. Boehm, I. Antonov, S. Baunsgaard, M. Dokter, R. Ginthör, K. Innerebner, F. Klezin, S. N. Lindstaedt, A. Phani, B. Rath, B. Reinwald, S. Siddiqui, S. Benjamin Wrede: SystemDS: A Declarative Machine Learning System for the End-to-End Data Science Lifecycle. **CIDR 2020**]

Data Cleaning Pipelines

- Automatic Generation of Cleaning Pipelines
 - Library of robust, parameterized data cleaning primitives (physical/logical)
 - Enumeration of DAGs of primitives & hyper-parameter optimization (HB, BO)



University	Country	1	Univer
TU Graz	Austria		TU Gra
TU Graz	Austria]	TU Gra
TU Graz	Germany]	TU Gra
IIT	India		IIT
IIT	IIT		IIT
IIT	Pakistan		IIT
IIT	India		IIT
SIBA	Pakistan		SIBA
SIBA	null		SIBA
SIBA	null	1	SIBA
0		d	



After imputeFD(0.5)

A	В	С	D	
0.77	0.80	1	1	
0.96	0.12	1	1	
0.66	0.09	null	1	
0.23	0.04	17	1	
0.91	0.02	17	null	
0.21	0.38	17	1	
0.31	null	17	1	
0.75	0.21	20	1	
null	null	20	1	
0.19	0.61	20	1	
0.64	0.31	20	1	

Dirty Data

A	В	C	D
0.77	0.80	1	1
0.96	0.12	1	1
0.66	0.09	17	1
0.23	0.04	17	1
0.91	0.02	17	1
0.21	0.38	17	1
0.31	0.29	17	1
0.75	0.21	20	1
0.41	0.24	20	1
0.19	0.61	20	1
0.64	0.31	20	1
	A 0.77 0.96 0.66 0.23 0.91 0.21 0.31 0.75 0.41 0.19 0.64	A B 0.77 0.80 0.96 0.12 0.66 0.09 0.23 0.04 0.91 0.02 0.21 0.38 0.31 0.29 0.75 0.21 0.41 0.24 0.19 0.61	A B C 0.77 0.80 1 0.96 0.12 1 0.66 0.09 17 0.23 0.04 17 0.91 0.02 17 0.21 0.38 17 0.31 0.29 17 0.75 0.21 20 0.41 0.24 20 0.19 0.61 20

After MICE

Multi-Level Lineage Tracing & Reuse



- Lineage as Key Enabling Technique
 - Trace lineage of operations (incl. non-determinism), dedup for loops/functions
 - Model versioning, data reuse, incremental maintenance, autodiff, debugging
- Full Reuse of Intermediates
 - Before executing instruction, probe output lineage in cache Map<Lineage, MatrixBlock>
 - Cost-based/heuristic caching and eviction decisions (compiler-assisted)

Partial Reuse of Intermediates

- Problem: Often partial result overlap
- Reuse partial results via dedicated rewrites (compensation plans)
- Example: stepIm

for(i in 1:numModels)
R[,i] = lm(X, y, lambda[i,], ...)

m_lmDS = function(...) {
 l = matrix(reg,ncol(X),1)
 A = t(X) %*% X + diag(1)
 b = t(X) %*% y
 beta = solve(A, b) ...}

es It overlap dicated ns) m>>n

t(X)

Х

ExDRa Project Overview





- Federated Backend
 - Federated data (matrices/frames) as meta data objects
 - Federated linear algebra, (and federated parameter server)
 - X = federated(addresses=list(node1, node2, node3), ranges=list(list(0,0), list(40K,70), ..., list(80K,0), list(100K,70)));

ex_{ra}

DDAI



Federated Requests: READ, PUT, GET, EXEC_INST, EXEC_UDF, CLEAR





Model Debugging [SIGMOD'21]

- Problem: Is 85% Accuracy good?
 - Intuitive slice scoring function
 - Exact top-k slice finding
 - $|S| \ge \sigma \land sc(S) > 0$
 - $\alpha \in (0,1]$

Properties & Pruning

- Monotonicity of slice sizes, errors
- Upper bound sizes/errors/scores
 → pruning & termination

Linear-Algebra-based Slice Finding

- Recoded matrix X, error vector e
- Vectorized implementation in linear algebra (join & eval via sparse-sparse matrix multiply)
- Local and distributed task/data-parallel execution



[Credit: sliceline, Silicon Valley, HBO]



slice error

SC

slice size







Summary and Q&A

- Course Goals
 - #1 Major data integration architectures
 - #2 Key techniques for data integration and cleaning
 - #3 Methods for large-scale data storage and analysis
- Programming Projects
 - Unique project in Apache SystemDS (teams or individuals), or
 - Exercise on data engineering and ML pipeline

Next Lectures

- O2 Data Warehousing, ETL, and SQL/OLAP [Oct 15]
- 03 Message-oriented Middleware, EAI, and Replication [Oct 22]

