Advanced Data Management (CSCI 680/490)

Machine Learning in Databases

Dr. David Koop





Checking Computational Results in Systems



D. Koop, CSCI 680/490, Spring 2022



Northern Illinois University

NIU





Repeatability Results



Figure 11: Study result. Blue numbers represent papers that were excluded from consideration, green numbers papers that are weakly repeatable, red numbers papers that are non-weakly repeatable, and orange numbers represent papers that were excluded (due to our restriction of sending at most one email to each author).

| OK ^{≤30} OK [≤] 130 64 | >30 OK ^{Auth} 23 | |
|---|------------------------------|--|
| | Notation | Number of papers |
| \mathcal{J} | HW | excluded due to replication requiring special hardware |
| Build fails 9 | NC | excluded due to results not being backed by code |
| | EX | excluded due to overlapping author lists |
| | BC | where the results are backed by code |
| | Article | where code was found in the paper itself |
| | Web | where code was found through a Web search |
| | EM yes | where the author provides code after receiving an email message |
| | EM ^{no} | where the author responds to an email message saying code cannot be provided |
| | EMø | where the author does not respond to email requests within two months |
| | OK ^{≤30} | where code is available and we succeed in building the system in \leq 30 minutes |
| | OK >30 | where code is available and we succeed in building the system in >30 minutes |
| | OK ^{Auth} | where code is available and we fail to build, and the author says the code builds with reasonable effort |
| \mathbf{M}^{\emptyset} $\mathbf{E}\mathbf{M}^{\mathrm{no}}$ 30 146 | Fails | where code is available and we fail to build, and the author says the code may have problems building |
| | | |













Excuses for not sharing

- Versioning
- Available Soon
- No Intention to Share
- Personnel Issues
- Lost Code
- Academic Tradeoffs
- Industrial Lab Tradeoffs
- Obsolete HW/SW
- Controlled Usage
- Privacy/Security
- Design Issues

D. Koop, CSCI 680/490, Spring 2022





Northern Illinois University



Examining 'Reproducibility in Computer Science'

- Repeat the experiment in reproducibility!
- Differences from original
- Shows issues with trying to classify experiments

| F | ի | 1 | r |
|---|---|---|---|
| | Γ |) | 1 |

All Others Purported Not 27%

- ported Not Building; 6% ••••• sputed; Not Checked
- Purported Building; Disputed; 2% •• Not Checked
 - Conflicting Checks! 0%
 - Misclassified 1% •
 - Purported Not Building But 14% ••••••••• Found Building
- Purported Building But Found 0% Not Building
 - Purported Not Building; 0% Confirmed
- Purported Building; Confirmed 0% •









Reproducible Research

- Science is verified by replicating work independently
- Replication Issues:
 - Requires many resources to replicate (Sloan Digital Sky Survey) - Requires significant computing power (Climate Model Simulation) - Requires too much time or very specific circumstances (Environment

 - Epidemiology)
- Reproducibility
 - Replication of the analysis based on the collected data (not replicating the data collection itself)
 - Better if we have the actual code or available executables _









Reproducibility Spectrum



D. Koop, CSCI 680/490, Spring 2022









7

10 Rules for Reproducible Computational Research

- Rule 1: For Every Result, Keep Track of How It Was Produced
- Rule 2: Avoid Manual Data Manipulation Steps
- Rule 3: Archive the Exact Versions of All External Programs Used
- Rule 4: Version Control All Custom Scripts
- Rule 5: Record All Intermediate Results, When Possible in Standardized Formats







10 Rules for Reproducible Computational Research

- Rule 6: For Analyses That Include Randomness, Note Underlying Random Seeds
- Rule 7: Always Store Raw Data behind Plots
- Rule 8: Generate Hierarchical Analysis Output, Allowing Layers of Increasing Detail to Be Inspected
- Rule 9: Connect Textual Statements to Underlying Results • Rule 10: Provide Public Access to Scripts, Runs, and Results











(Database) Reproducibility Research Topics

- Design and Management of Experiment Repositories
- Querying and Searching Experiments
- Mining Experiments







Notebook Reproducibility

- Use notebooks from Github (~1 million) - Unambiguous cell order? 81.99%
- Study notebook dependencies
 - Dependencies Available? 13.72%
 - Dependencies Install? 5.03%
- Study notebook executability
 - Execute: 24.11% of unambiguous cell order
 - Matched results: 4.03%







Dataflow Notebooks: Resolve Notebook Ambiguities

| | | | | | | | | | | 1 | | | | | | | |
|----------------|--|---|------|-------------|-------|-----------|---------------|--------------|---|---|--------|----------------|-------|-------------|-----|-------------|---|
| In [d51f8eab]: | <pre>import pandas as pd df = pd.read_csv('guardian-top100-female-2019.csv')</pre> | | | | | | | In [over30]: | <pre>df = df\$full[df\$full.Age >= 31]</pre> | | | | | | | | |
| df: | N | ame F | Rank | Position | Age o | n 1 Dec 2 | 019 1 | Nationality | | df: | | Name | Rank | Position | Age | Nationality | _ |
| | 0 Sam | Kerr | 1 | Forward | | | 26 | Australia | | | 2 1 | Megan Rapinoe | 3 | Midfielder | 34 | USA | |
| | | | | | | | | | | | | | | | | | |
| | 99 Luc | Imila | 100 | Forward | | | 25 | Brazil | | | 96 | Cláudia Neto | 97 | Midfielder | 31 | Portugal | |
| | 100 rows × 5 columns | | | | | | | | 19 rows × 5 columns | | | | | | | | |
| In [full]: | df = d | <pre>df = df.rename(columns={'Age on 1 Dec 2019': 'Age'})</pre> | | | | | | | In [under25]: df: | <pre>df = df\$full[df\$full.Age <= 24]</pre> | | | | | | | |
| df. | N | Name Bank Position Age Nationality | | | | | | | | Name | Rank | Position | Age I | Nationality | | | |
| ur. | 0 Sam | Kerr | 1 | Forward | 26 | Australi | <u>,</u> а | | | | 3 / | Ada Hegerberg | 4 | Forward | 24 | Norway | |
| | | | • | | | | | | | | | | | | | | |
| | 99 Luc | lmila | 100 | Forward | 25 | Braz | il | | | | 98 | Lena Oberdorf | 99 | Midfielder | 17 | Germany | |
| | 100 rows > | 5 colur | mns | | 20 | 2:42 | | | | | 25 rov | ws × 5 columns | | | | | |







Final Exam

- Monday, May 9, 4:00-5:50pm, PM 153
- Similar format
- More comprehensive (questions from topics covered in Test 1 & 2)
- Will also have questions from graph/spatial/temporal data, provenance, reproducibility, machine learning





Improving Databases





LEARNED AND **SELF-DESIGNING** DATA STRUCTURES Data Systems and AI Lab Stratos Idreos & Tim Kraska



Algorithms rely on the order of data





D. Koop, CSCI 680/490, Spring 2022

ALGORITHMS

[7,4,2,6,1,3,9,10,5,8]











Data systems rely on algorithms

DATA SYSTEMS ALGORITHMS

D. Koop, CSCI 680/490, Spring 2022











17

Data structures define performance



D. Koop, CSCI 680/490, Spring 2022



register = this room caches = this city memory = nearby city disk = Pluto

Jim Gray, Turing Award 1998







18

Database Questions

How do I make my **data system** run x times as fast?



How do I extend the **lifetime** of my hardware?

How to accelerate **statistics** computation for data science/ML?



D. Koop, CSCI 680/490, Spring 2022



How do I minimize my **bill** in the **cloud**?





How do I train my **neural network** x times faster?







Tradeoffs in each structure









New Applications Demand Change

existing systems need to change too

















"Traditional" Database Research











Self-designing systems







D. Koop, CSCI 680/490, Spring 2022

without coding or accessing the h/w

algorithms



performance











SageDB: a learned database system

G. Leclerc, S. Madden, H. Mao, and V. Nathan

D. Koop, CSCI 680/490, Spring 2022

T. Kraska, M. Alizadeh, A. Beutel, E. H. Chi, J. Ding, A. Kristo,





Learned Data Structures and Algorithms









Discussion

- Is this the future?
- What about comparison baselines?
- Lots of work being done in this area











Benchmarking Learned Indexes



D. Koop, CSCI 680/490, Spring 2022

[<u>R. Marcus et al.</u>, 2021]









Multi-Dimensional Indexing











Query Optimization











Reminders

- Final Exam Review Wednesday (come with questions!)
- Final Exam on Monday, May 9 from 4-5:50pm





