

CSE P 590

Building Data Analysis Pipelines

Fall 2024

Big data



Today

- More on big data processing
- Wrap up and open discussion
- In-class 5

More on big data processing

Recap: Big data

Characteristics

- *Volume*: data sets are (too) big → distributed analysis
 - *Variety*: data formats: structured, semi-structured, unstructured
 - *Velocity*: data changes rapidly → real-time analysis
-
- *Variability*: meaning of data changes
 - *Veracity*: noisy data (tradeoff between noisy and useful)
 - *Value*: informed decision making (value of collected data)

Recap: Big data challenges

Compute bound

- Compute-intensive simulations
- Real-time processing
- High-volume data

Memory bound

- Data exceeds memory on a single machine

I/O bound

- Different data sources (structured, semi-structured, unstructured)
- Data sharing among different processes

Recap: sparklyr (a Spark DSL)

Distributed data storage

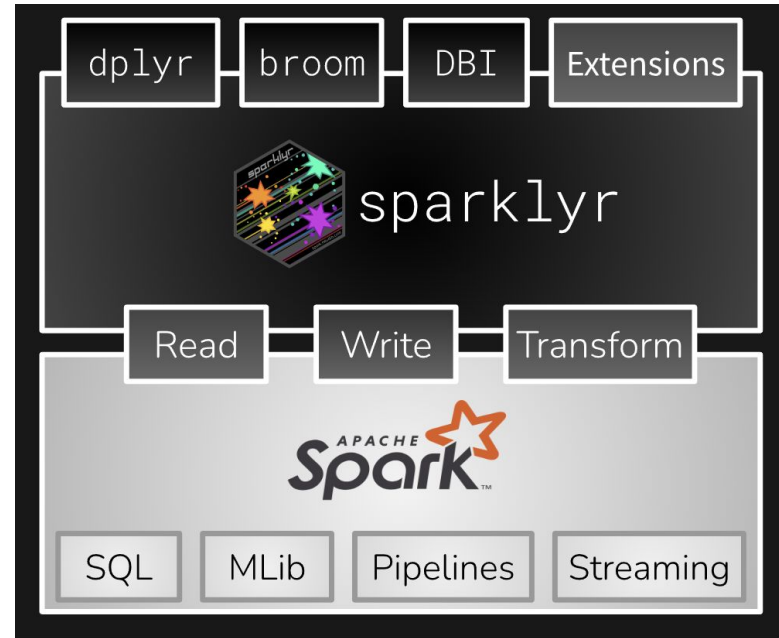
- A dplyr backend
- Supports SQL (like dbplyr)
- Lazy evaluation

Distributed data processing

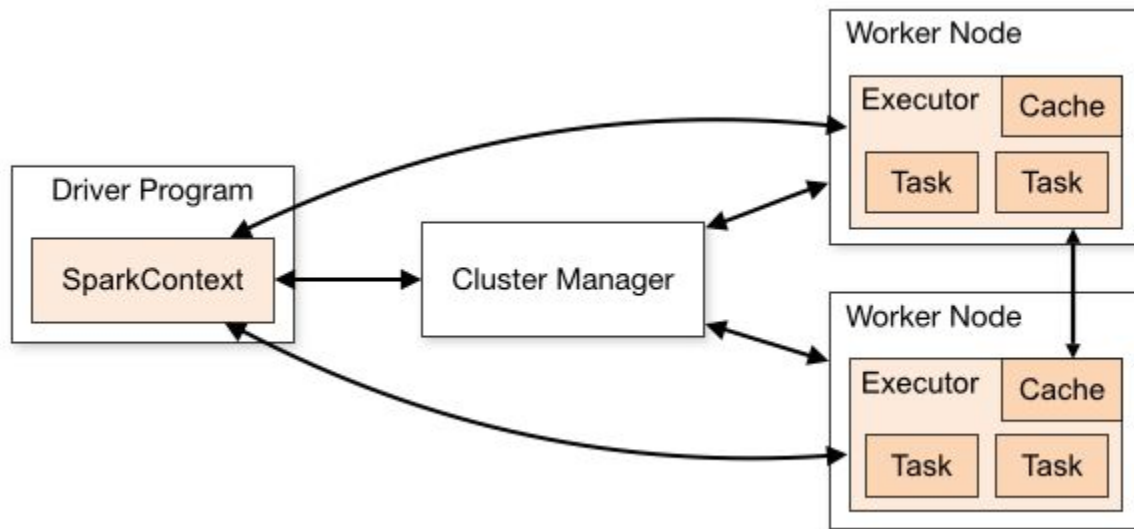
- Support for many data formats

Modeling/ML

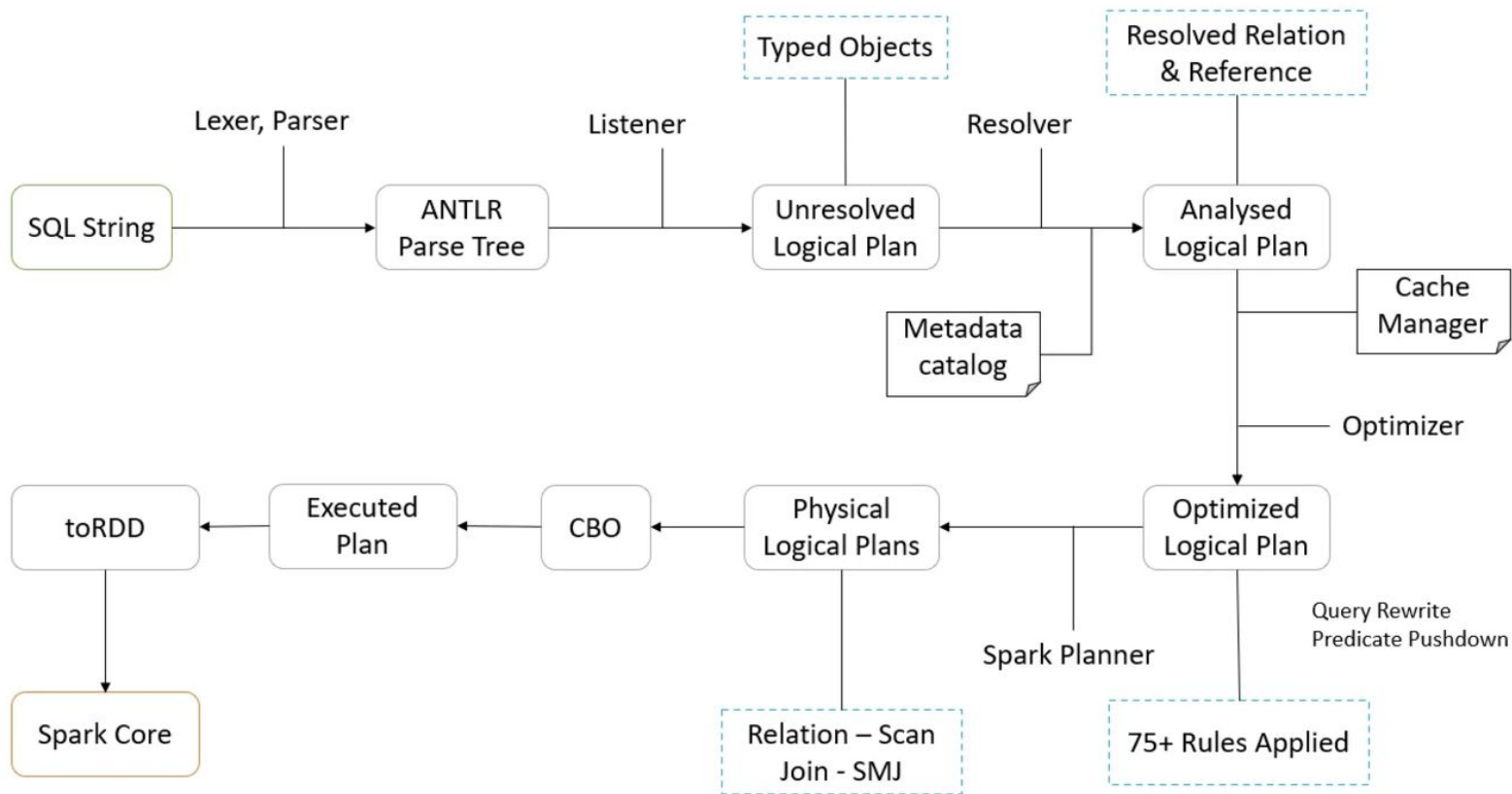
- Support for many common model types



Spark: architecture



Spark: query optimization



Dplyr backends: Spark vs. DB

- Sparklyr and dbplyr as dplyr backends
- Automatic translation from dplyr verbs to SQL queries
- Lazy evaluation to minimize network traffic

(If your data fits in memory there is no advantage to putting it in a database: it will only be slower and more frustrating.)

Most of the time you don't need to know anything about SQL, and you can continue to use the dplyr verbs that you're already familiar with:

However, in the long-run, I highly recommend you at least learn the basics of SQL. It's a valuable skill for any data scientist, and it will help you debug problems if you run into problems with dplyr's automatic translation.

<https://dbplyr.tidyverse.org/articles/dbplyr.html>

<https://spark.posit.co/guides/dplyr.html>

Live demo: Spark vs. DB

Open discussion

Why R for this course?



1. Leveling the playing field (Python knowledge distribution is bimodal)
2. (My opinion) dplyr, ggplot2, etc. are superior: great for data wrangling/viz.
3. Focusing on (non-ML) analysis concepts
4. Many specialized data science packages in R

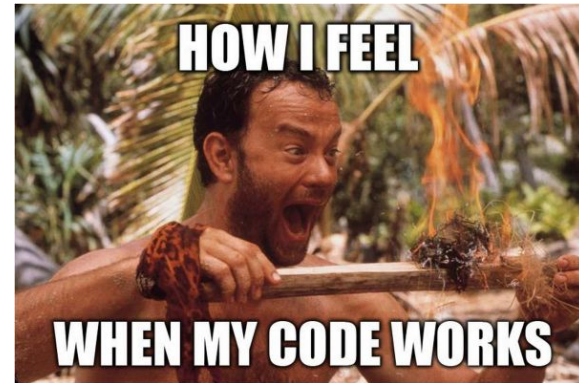
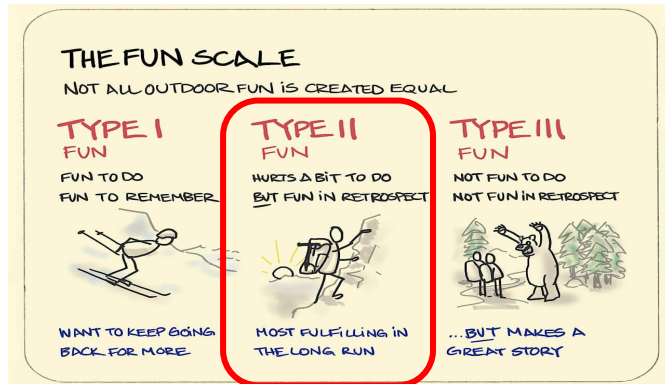
Pick the right tool/ecosystem for the job!

Open discussion

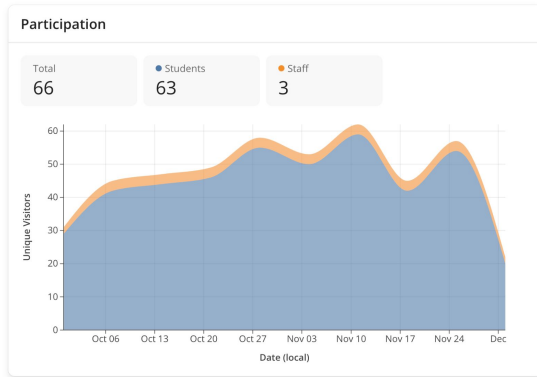
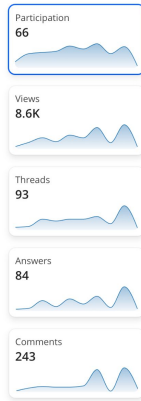
Example topics

- Highlights and frustrations
- Open questions and confusions
- Experiences with AI-assisted coding
- R vs. Python (your conclusions)

10 weeks flew by!



Thanks for a great quarter!



Document tabs +

Tab 1

Building Data Analysis Pipelines - Course Notes

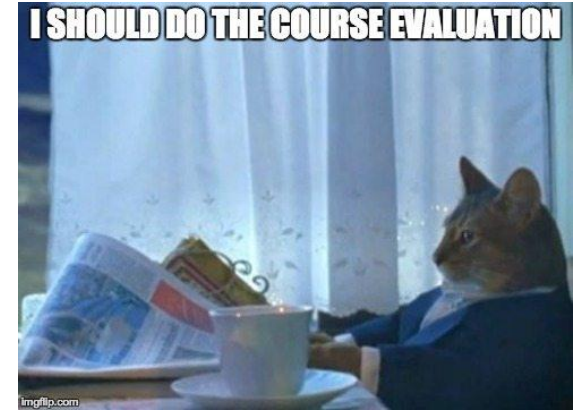
Canvas notes

Canvas issues this quarter

- Canvas has been a huge pain this quarter.
- I don't recall similar issues being nearly as bad in the past.
- Group sets are not visible to students (no individual tabs for group sets). Instead, students only see a list of all consolidated groups - for our course this means 5 x 50 groups. Students need to scroll all the way down to find a group they can self-assign themselves to for an exercise.

Canvas groups

| Group | Members | Group Set | Group Set | Group Set | Group Set | Group Set |
|------------------------------|---------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|
| 10:07: Validity | 50 | 10:07: Validity | 10:07: Validity | 10:07: Validity | 10:07: Validity | 10:07: Validity |
| 10:14: Data wrangling | 50 | 10:14: Data wrangling | 10:14: Data wrangling | 10:14: Data wrangling | 10:14: Data wrangling | 10:14: Data wrangling |
| 10:21: Statistical model... | 50 | 10:21: Statistical model... | 10:21: Statistical model... | 10:21: Statistical model... | 10:21: Statistical model... | 10:21: Statistical model... |
| 10:28: Significance and... | 50 | 10:28: Significance and... | 10:28: Significance and... | 10:28: Significance and... | 10:28: Significance and... | 10:28: Significance and... |
| 11:04: Advanced statist... | 50 | 11:04: Advanced statist... | 11:04: Advanced statist... | 11:04: Advanced statist... | 11:04: Advanced statist... | 11:04: Advanced statist... |
| 11:18: Data visualization... | 50 | 11:18: Data visualization... | 11:18: Data visualization... | 11:18: Data visualization... | 11:18: Data visualization... | 11:18: Data visualization... |
| 11:25: Big data (1) | 50 | 11:25: Big data (1) | 11:25: Big data (1) | 11:25: Big data (1) | 11:25: Big data (1) | 11:25: Big data (1) |
| 12:02: Big data (2) | 50 | 12:02: Big data (2) | 12:02: Big data (2) | 12:02: Big data (2) | 12:02: Big data (2) | 12:02: Big data (2) |



<https://uw.iasystem.org/survey/297879>

In-class 5