

Lab #1: Data Wrangling and Exploration

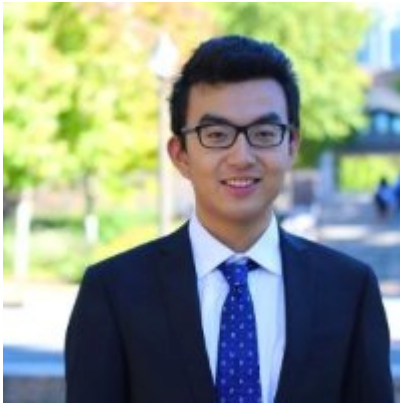
Everything Data
CompSci 216 Spring 2015



DUKE
COMPUTER SCIENCE

Announcements (Wed. Jan. 14)

- Welcome our UTAs!



Billy Wan



Eric Wu

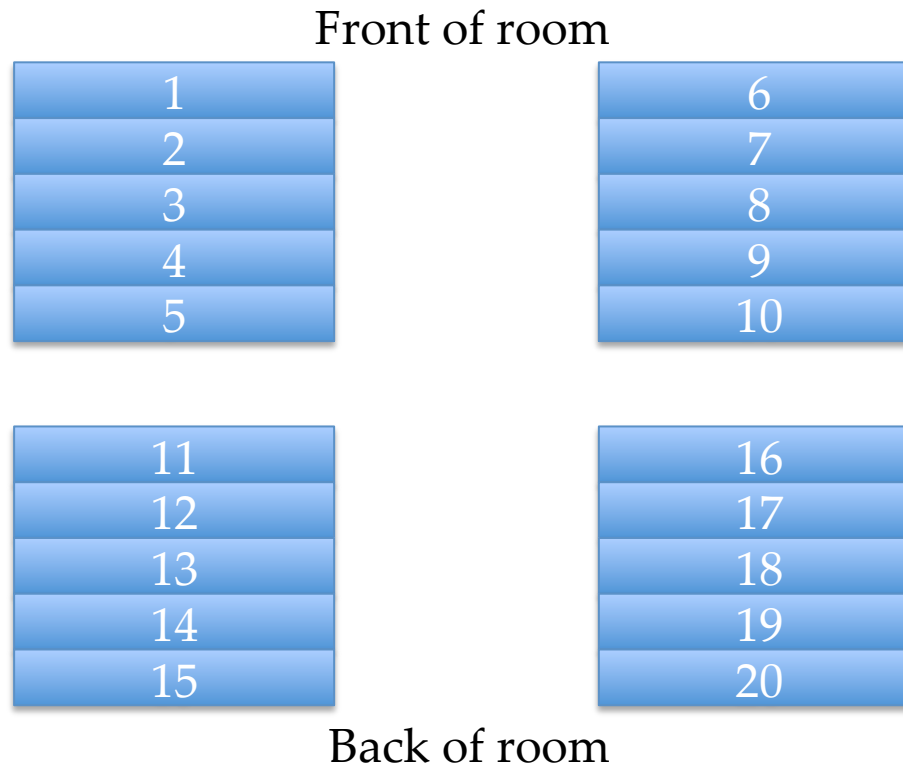


Kevin Wu

- Most of our office hours have been posted on the website

Seat assignment

See course website (under “Schedule”) for team assignment



Format of this lab

- Discuss last homework
 - Continue working on your setup if needed
- Work on new exercises
 - Raise your hand to get your answer checked by the course staff
 - Win challenge to get little prizes!
- Summarize lessons learned (10 min.)

Lesson learned: reality check

- Data is messy
- Reality is messy
- Garbage in, garbage out



☞ *One reason why you need data wrangling*

Lesson learned: “abstraction”

- More structure and semantics \Rightarrow more powerful analysis, and
- Different data models \Rightarrow different questions and processing methods



☞ Another reason you need data wrangling

- Examples in this lab
 - Character strings: regular expressions
 - Dates: meaningful range conditions, subtraction
 - Tables: filtering, grouping, counting, sorting...

Lesson learned: UI or not UI

- Exploratory analysis goes a long way
 - You can get fair amount of insight without a single line of code



- Interactivity is nice
- Easy to learn,
difficult to master
 - Tool-specific recipes vs.
universal primitives

What's next

- No class next Monday (MLK Day)
- No homework until next Wednesday (lecture on relational data processing)
 - Resolve any remaining setup issues
 - Try to gain some efficiency with VM and shell