# cloudera

# DATA SCIENTIST TRAINING

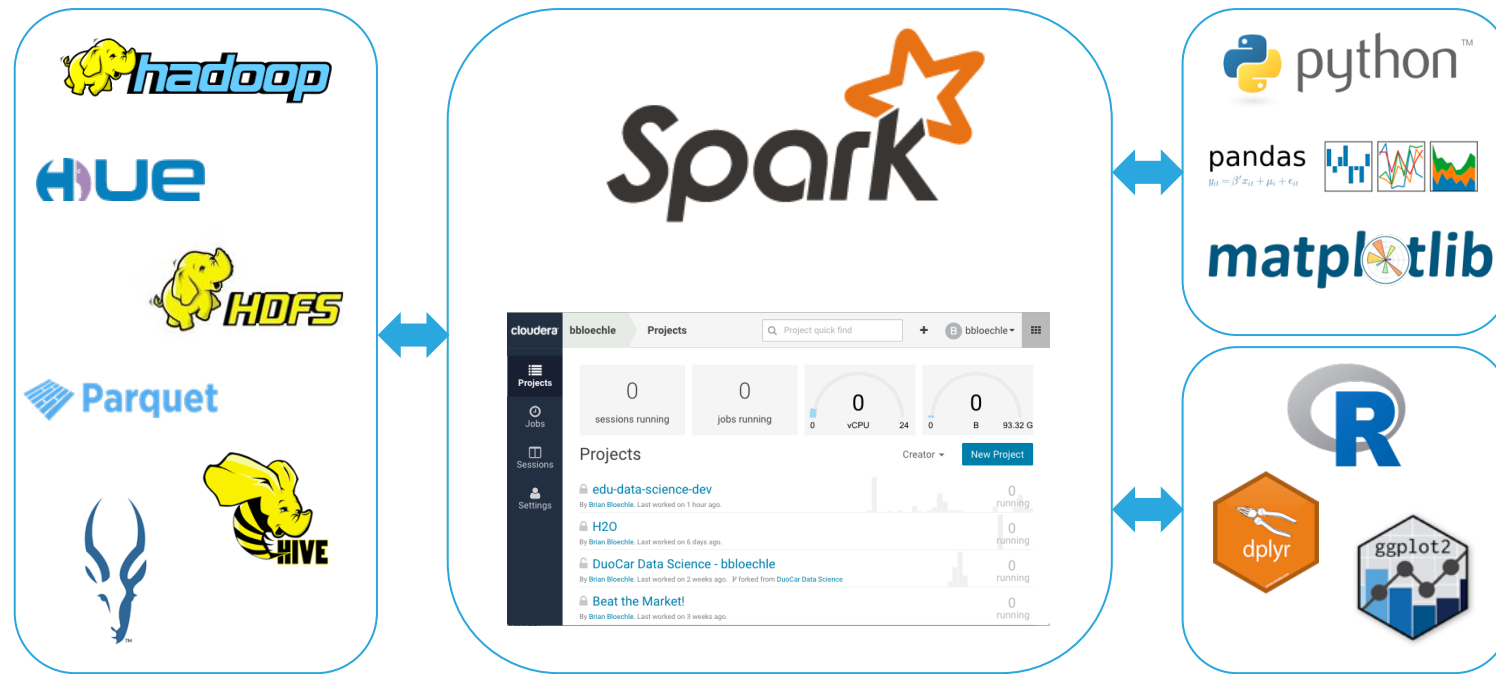Kai Voigt, Senior Instructor, Cloudera University

# HIGHLIGHTS

A course on Data Science and Machine Learning at Scale

- Four-Day workshop using interactive code modules instead of slides
- Designed for both
  - Data Scientists with limited Hadoop experience
  - Hadoop Data Engineers with limited Data Science experience
- Focused on industry-leading tools
  - Hadoop, Spark, Python, R, SQL
  - Cloudera Data Science Workbench (CDSW)

# OBJECTIVES

- Data Engineering
  - Reading and Writing Data
  - Inspecting Data Quality
  - Cleansing and Transforming Data
  - Combining and Splitting Data
  - Summarizing and Grouping Data
  - Exploring and Visualizing Data

- Data Science and Machine Learning
  - Extracting and Selecting Features
  - Building and Evaluation of
    - Regression Models
    - Classification Models
    - Clustering Models
  - Cross-Validating Models and Tuning Hyperparameters
  - Building Pipelines
  - Deploying Models
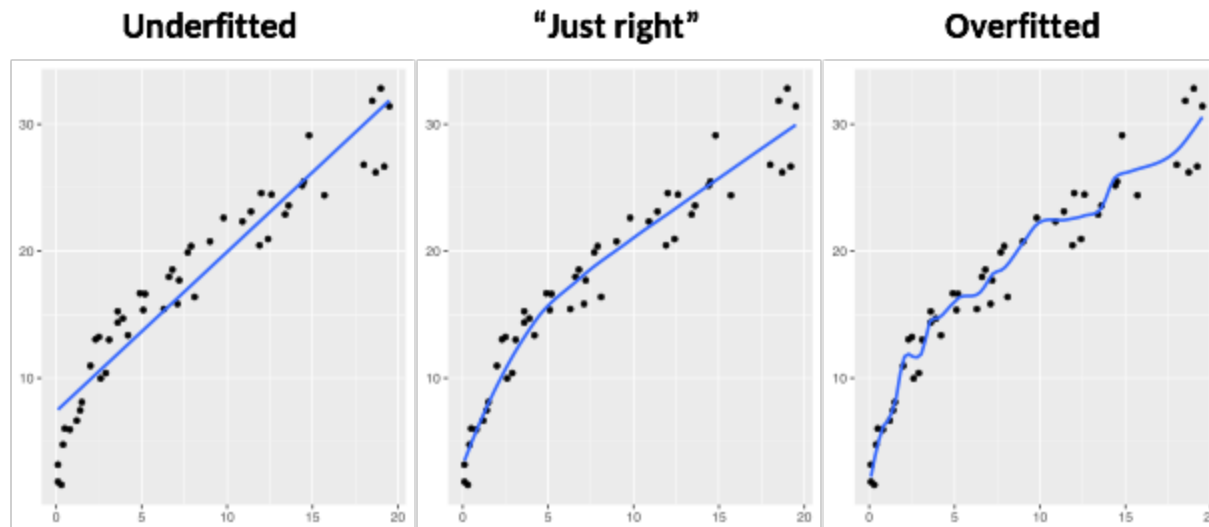
# ARCHITECTURE AND TOOLS
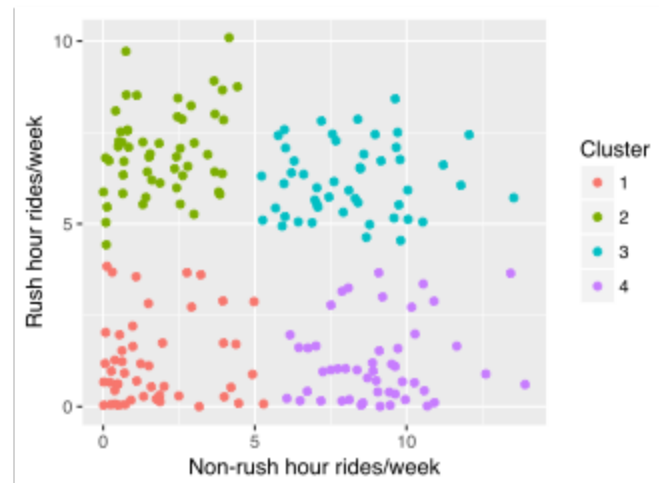
# OUR STORY

- Ridesharing company with lots of data
  - Drivers
  - Riders
  - Rides
  - Ride Routes
  - Reviews
  - Weather
  - Demographics

duocar

# DISCUSSION EXAMPLE

# ANOTHER EXAMPLE

# DEMO TIME

# THANK YOU