# DataLab: Introducing Software Engineering Thinking into Data Science Education at Scale

**Yang Zhang**, Tingjian Zhang, Yongzheng Jia, Jiao Sun, Fangzhou Xu, and Wei Xu

Institute of Interdisciplinary Information Sciences, Tsinghua University

Department of Computer Science and Technology, Shandong University

# Overview - Backgrounds

▶ Data science

Data scientist became the best job in the US in 2016

▶ Data Science Education

Ubiquitous in Universities and Online Education



50 Best Jobs in America

This report ranks jobs according to each job's Glassdoor Job Score, determined by combining three factors: number of job openings, salary, and overall job satisfaction rating.

Employers: Want to recruit better in 2017? Find out how.

United States ▾     2017 ▾                    12k Shares

1   Data Scientist

4.8 / 5          4.4 / 5
Job Score      Job Satisfaction

$110,000       4,184
Median Base Salary   Job Openings

View Jobs

[1]: 25 Best Jobs in America. https://www.glassdoor.com/List/Best-Jobs-in-America-LST KQ0,20.htm

# Overview - Challenges

Students

Instructors

- Lack formal computer science training
- Hard to set up coding tools
- Confused with data/code versions

- Time-consuming to setup tools

- Hard to scale teaching methodologies
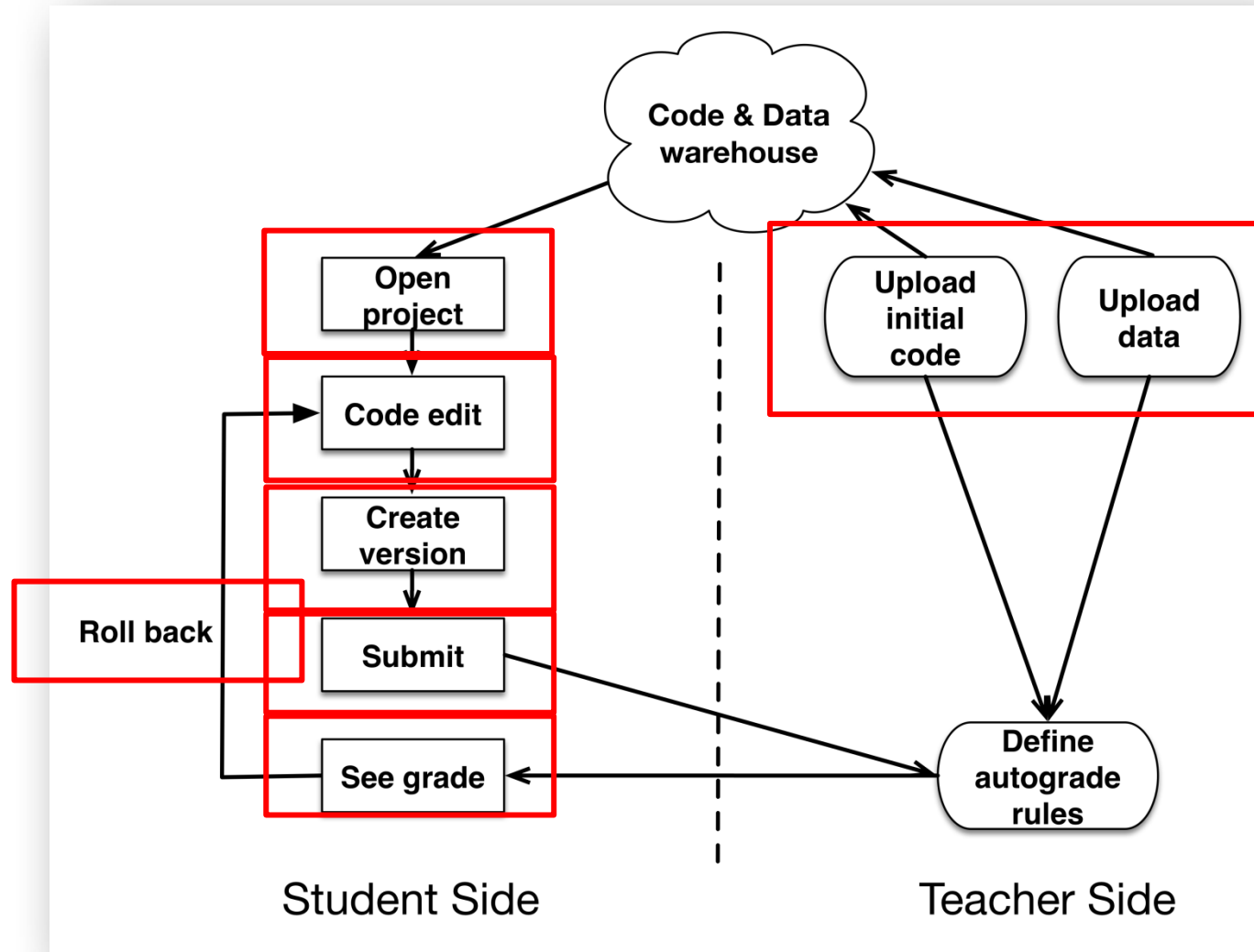
# Differences between a DS and SE project

# Our solution

- DataLab
  - Integrates code, data and execution management into a single system
  - Creates links among code, data, parameters and their revisions
  - Provides a scalable system
  - Allows students to share their code, data, results with any versions.

# Easy to set up a project

# A project summary page

▶ Data

▶ Code

▶ Project push commit

# Separate config and parameters from code

# Online development environment

# Creating code/data versions and autograding

# Version management

# Team collaboration

# Instructor tools

# DataLab is scalable

- Data management system

- Scalable execution environment

- Extensible APIs

# Evaluation - Deployment

- DataLab: 3 machines
  - 8 cores
  - 16 GB memory
  - 80GB of hard disk storage



[1]: Kaggle. https://www.kaggle.com

# Evaluation : in-classroom experiment

▶ A graduate-level introductory data science course with **81** students and **20** volunteers

▶ Classical Kaggle[1] competition project: *Titanic Machine Learning from Disaster*

  ▶ Predict survivors from gender, age, cabin class, and other information

  ▶ **1,979** different versions of code submissions

[1]: Kaggle. https://www.kaggle.com

# Log analysis



Fig 1. Relation between number of submissions and accuracy



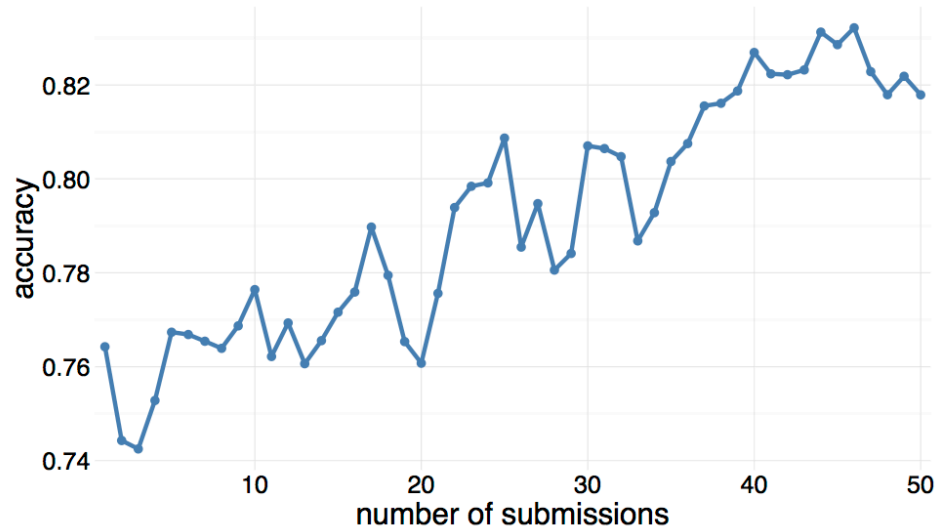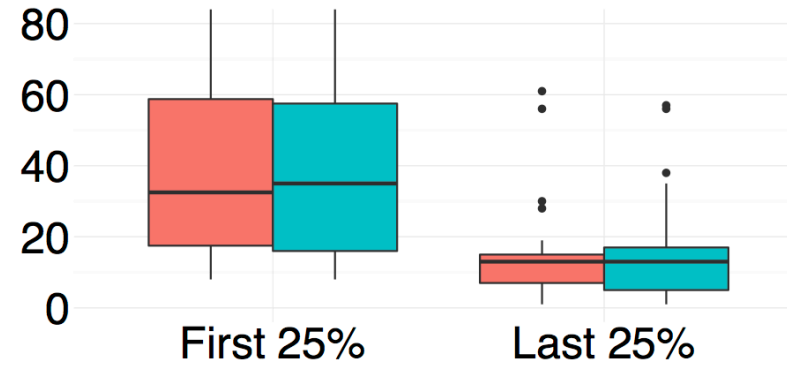Fig 2. How many times did students push and submit their code given their ranks?
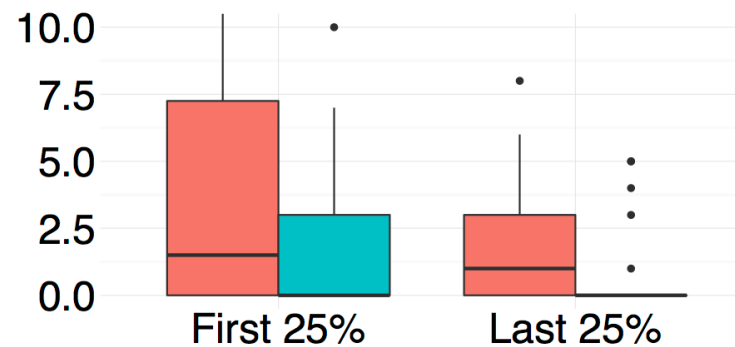


Fig 3. How many times did students check branches and reset their code given their ranks?

# Survey results

- 18 subjective questions

- The survey has 3 parts
  - Students' coding experience
  - Students' opinions
  - Students suggestions



Fig 1. Is DataLab helpful for learning data analysis techniques?

- **92** out of **101** students indicate that they will continue to use DataLab for their future data science projects

# Conclusion

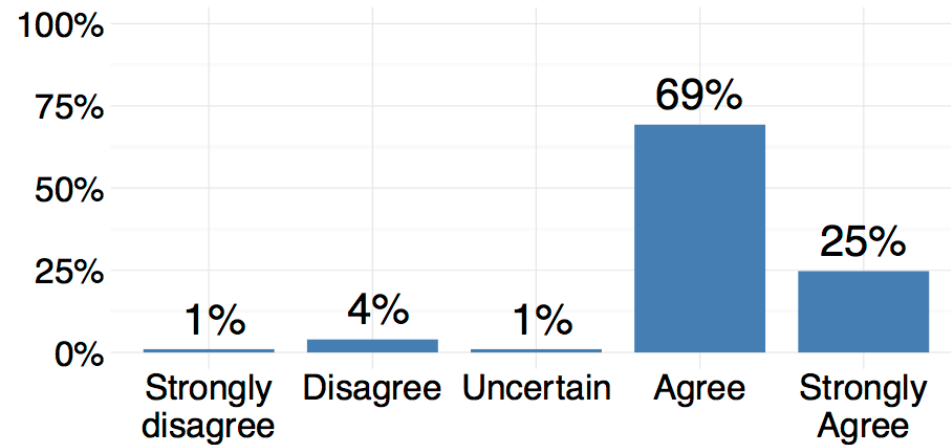## Datalab: introducing SE Thinking to DS Education

Save instructors' time

Improve students' development efficiency

Manage data/code/executi on automatically

Can scale at low cost