# Integrating R & Python into a Data Science program
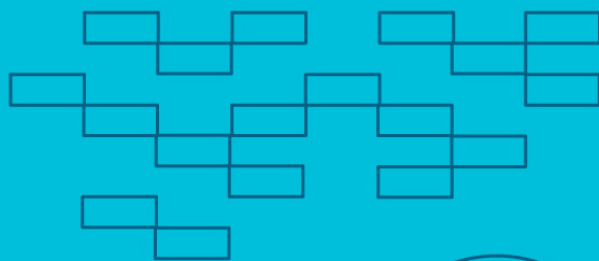
Tiffany Timbers & Ian Flores Siaca

University of British Columbia

# UBC Master of Data Science

## FALL
### SEP - DEC

**Block 1** (4 weeks)

- 511 - Programming for Data Science
- 521 - Computing Platforms for Data Science
- 542 - Communication and Argumentation
- 551 - Descriptive Statistics and Probability for Data Science

**Block 2** (4 weeks)

- 523 - Data Wrangling
- 531 - Data Visualization I
- 512 - Algorithms and Data Structures
- 552 - Statistical Inference and Computation I

**Block 3** (4 weeks)

- 561 - Regression I
- 532 - Data Visualization II
- 571 - Supervised Learning I
- 513 - Databases and Data Retrieval

## WINTER
### JAN - APR

**Block 4** (4 weeks)

- 562 - Regression II
- 573 - Feature and Model Selection
- 572 - Supervised Learning II
- 522 - Data Science Workflows

**Block 5** (4 weeks)

- 563 - Unsupervised Learning
- 553 - Statistical Inference and Computation II
- 524 - Collaborative Software Development
- 574 - Spatial and Temporal Models

**Block 6** (4 weeks)

- 575 - Advanced Machine Learning
- 541 - Privacy, Ethics and Security
- 554 - Experimentation and Causal Inference
- 525 - Web and Cloud Computing

## SPRING
### MAY - JUN

CAPSTONE
PROJECT

(8 weeks)
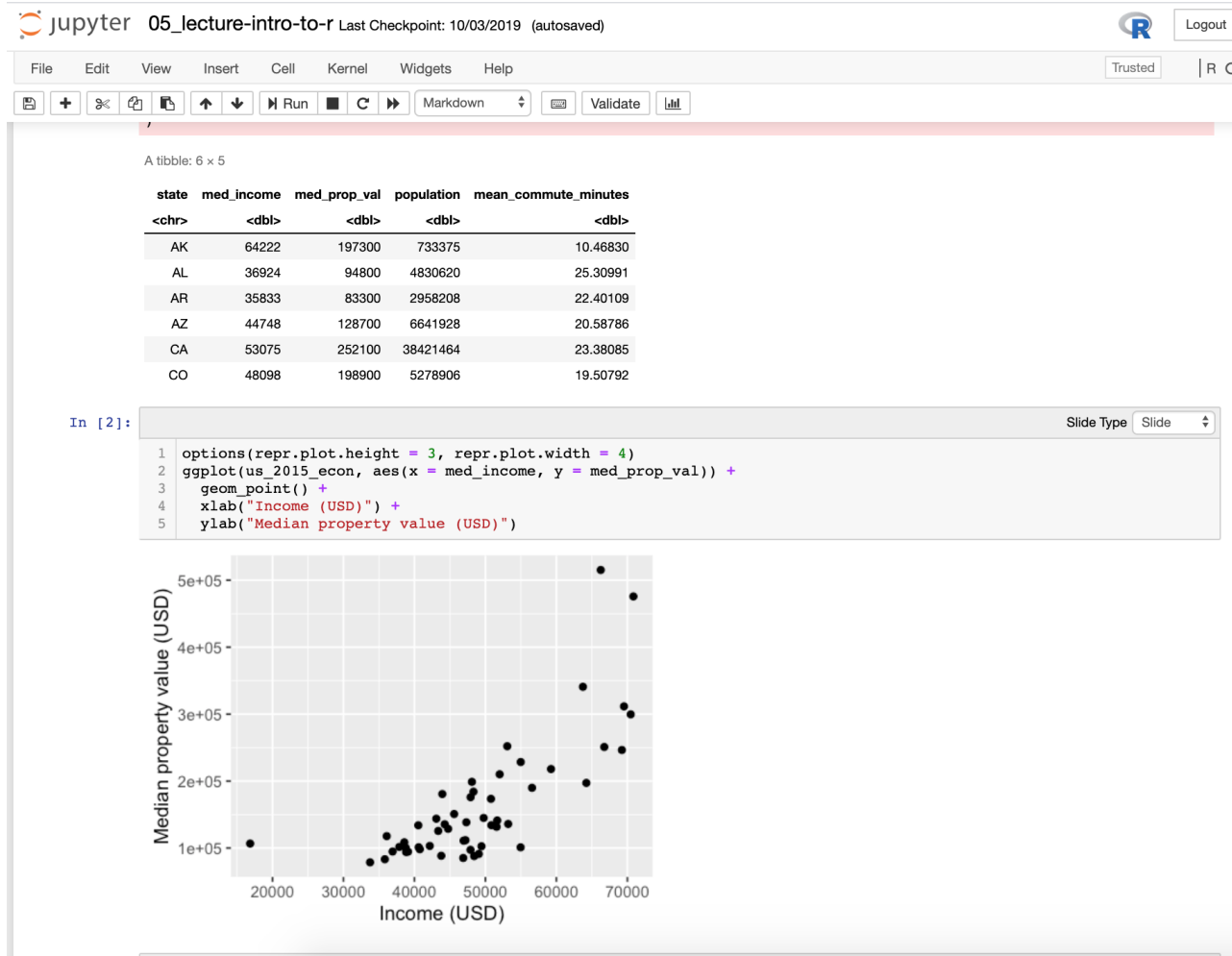
Languages used: R, Python, R & Python

# UBC Master of Data Science program

## Tools we teach for R & Python harmony

- RStudio
- Jupyter
- knitr & R Markdown
- feather file format
- reticulate
- Make
- Docker
- plotly Dash



reticulate

GNU Make

# Example 1: R in Jupyter!

# Example 2: RStudio as a Python IDE!

# Example 3: GNU Make for polyglot automation!

```
 ⚙ Makefile ×
 ◀▶ ▾ | 📁 | 💾 | 🔍

  1▾  # Tiffany Timbers, Nov 2018
  2   # usage: make all
  3
  4   # run all analysis
  5   all: doc/count_report.md
  6
  7   # make dat
  8▾  results/isles.dat: data/isles.txt src/wordcount.py
  9     python src/wordcount.py data/isles.txt results/isles.dat
 10▾  results/abyss.dat: data/abyss.txt src/wordcount.py
 11     python src/wordcount.py data/abyss.txt results/abyss.dat
 12▾  results/last.dat: data/last.txt src/wordcount.py
 13     python src/wordcount.py data/last.txt results/last.dat
 14▾  results/sierra.dat: data/sierra.txt src/wordcount.py
 15     python src/wordcount.py data/sierra.txt results/sierra.dat
 16
 17   #create plot
 18▾  results/figure/isles.png: results/isles.dat src/plotcount.py
 19     python src/plotcount.py results/isles.dat results/figure/isles.png
 20▾  results/figure/abyss.png: results/abyss.dat src/plotcount.py
 21     python src/plotcount.py results/abyss.dat results/figure/abyss.png
 22▾  results/figure/last.png: results/last.dat src/plotcount.py
 23     python src/plotcount.py results/last.dat results/figure/last.png
 24▾  results/figure/sierra.png: results/sierra.dat src/plotcount.py
 25     python src/plotcount.py results/sierra.dat results/figure/sierra.png
 26
 27   # make count_report
 28▾  doc/count_report.md: doc/count_report.Rmd results/figure/isles.png results/figure/abyss.png results/figure/last.png
 29     Rscript -e "rmarkdown::render('doc/count_report.Rmd')"
 30
 31   #Clean up intermediate files
 32▾  clean:
 33     rm -f results/isles.dat
 34     rm -f results/abyss.dat
 35     rm -f results/last.dat
 36     rm -f results/sierra.dat
 37     rm -f results/figure/isles.png
 38     rm -f results/figure/abyss.png
 39     rm -f results/figure/last.png
 40     rm -f results/figure/sierra.png
```
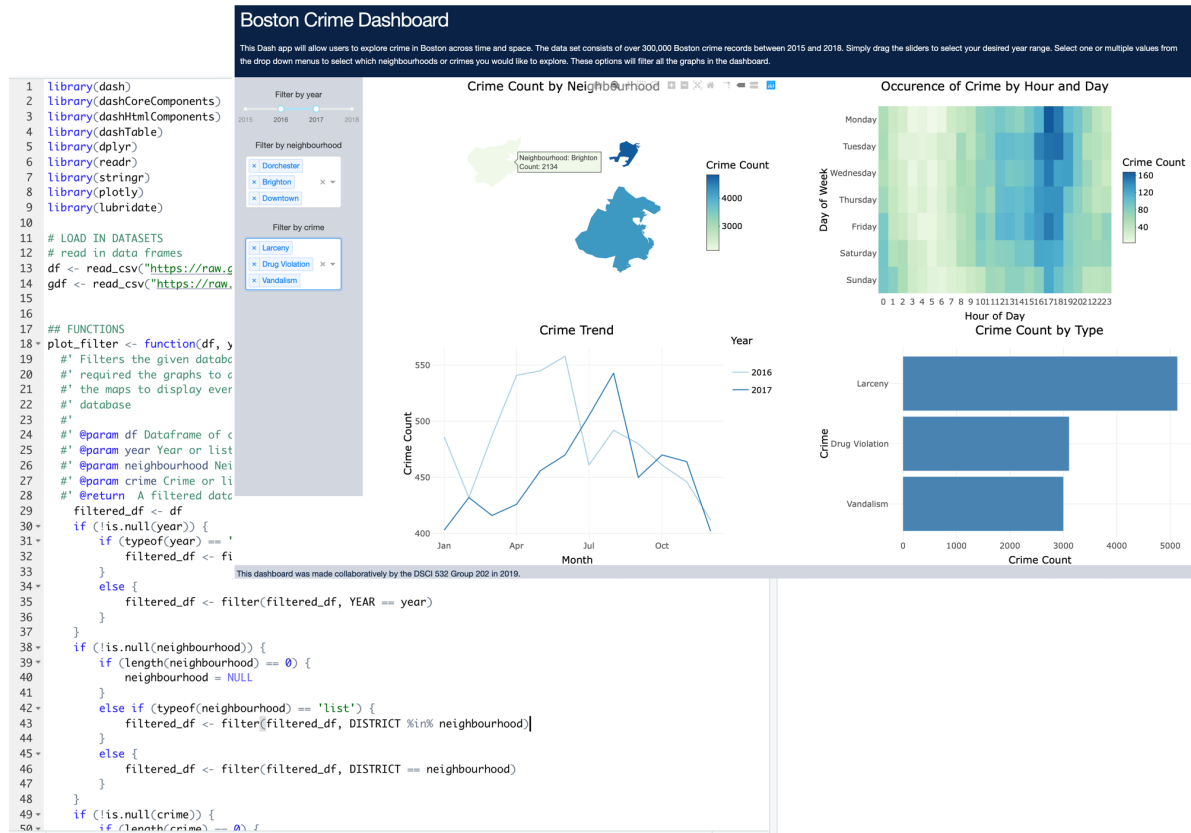
# Example 4: Docker for polyglot reproducibility!

# Example 5: Dashboards in R with plotly Dash!

# Pedagogical challenges (and solutions!) for teaching both R & Python

**Problem 1:** Mixed proficiencies of previous R & Python programming skills between students

**Solution 1:** Optional questions to challenge more advanced students, and extra practice questions with feedback to support novices.

# Pedagogical challenges (and solutions!) for teaching both R & Python

**Problem 2:** Dual task interference

**Solution 2:** Learning outcomes in the program include comparing and contrasting the diffences between the languages (i.e., we spend a lot of time teaching and assessing whether the students know this).

# Pedagogical challenges (and solutions!) for teaching both R & Python

**Problem 3:** Memory decay during breaks in practice

**Solution 3:** All blocks in the program have courses that require students to use R & Python

# Take homes:

Tips for integrating R & Python into a Data Science program:

- Carefully choose tools that work well with both languages, and skip the ones that don't.

- Expect students to have a heterogeneous knowledge base that may differ between languages, and design exercises to address this.

- Teach the R'isms and the Python'isms and have the students compare and contrast them. Also, asses them on this!

- Structure the program so students repeatedly practice both languages, avoid gaps in one language if possible!

# Thanks!

UBC MDS public resources: https://github.com/UBC-MDS/public