# Using Cloud Computing Resources for Big Data and Code Reproducibility

## —

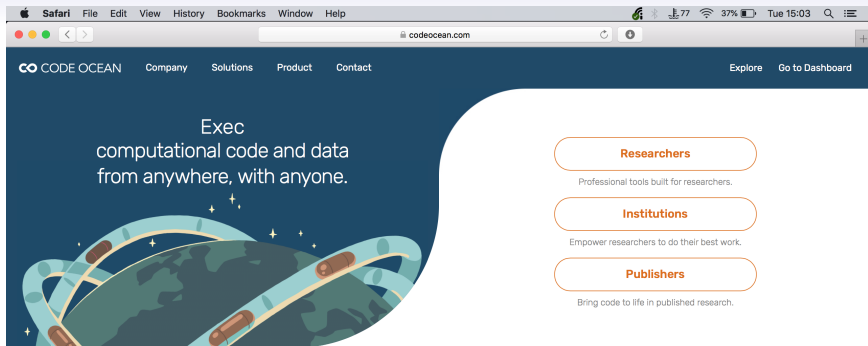## Code Ocean

**Winter Institute in Data Science and Big Data**
Simon Heuberger
9 January 2020

AMERICAN UNIVERSITY
W A S H I N G T O N , D C

Code Ocean is a research collaboration platform.

We support researchers from the beginning of a project through publication. With direct access to cloud computing and reproducibility best practices built in, no extra software or hardware is needed.

# What is a capsule?

- Code + (optional) Data + Computational environment = The minimum required for computational reproducibility
- At the core of the capsule: Docker image
- Potential to revolutionize data reproduction and replication

# What is a capsule?

- Code?
- Data?
- Computational environment?

# Computational environment

- Base environment
  - Operating system
  - Programming language
- Packages/Dependencies

# Example: Actual capsule

# Exercises

On our GitHub repository: 08-03-cloud-authors_how_to.pdf

1. Create a Code Ocean account at https://codeocean.com
2. Create a new capsule
3. Add R as the base environment
4. Add packages along with their specific versions
5. Upload data
6. Upload one .R file

    - The file should read in the data
    - The file should create and save one plot with ggplot2

7. Create a run script for the uploaded .R file
8. Upload a very rudimentary Readme file
9. Commit the changes
10. Execute one successful reproducible run

# Extra exercises

- Add Python 3 to your R base environment via apt-get
- Include a full log of the files in run
- Save your plot to the subfolder /figures
- Create another new capsule by importing your Git repository