

How Usability and Reproducibility in Software Improves Teaching and Research

Soren Harner, RC2AI

SDSS 2021, Soren Harner, RC2AI

1

Reproducibility

Given the same raw data, can you follow the steps and understand the assumptions of how the authors arrived at their conclusion?

SDSS 2021, Soren Harner, RC2AI



Why does it matter?

- Improves trust and transparency in science and beyond
- Promotes learning, habits, and building on others' work

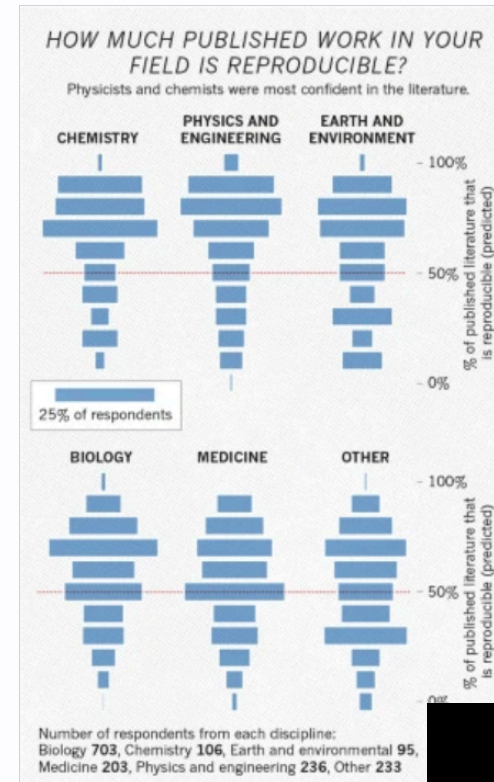
SDSS 2021, Soren Harner, RC2AI



Trust in research

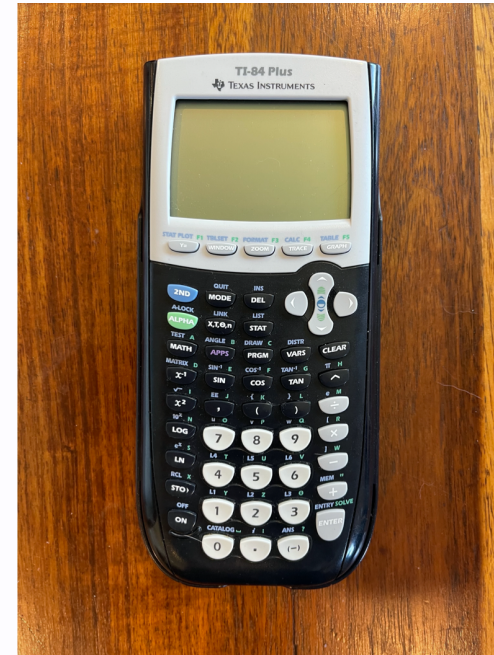
- [Nature 2016 Survey](#) 52% say there is crisis
- [Nature 2021 Survey](#) confidence in science leads to vaccination

SDSS 2021, Soren Harner, RC2AI



Learning by doing

- Data used in all fields
- Computational sciences
- Larger models and big data
- Need to automate



In 2021, really?



SDSS 2021, Soren Harner, RC2AI

Why is it hard?

- Not WYSIWYG
- Hacker mindset
- Accessible data, code
- Replicate environment
- Cloud operations
- Trust chain for security

Zooter.xlsx

Operations Analytics MOOC

Maximize $150R + 160N$
 subject to
 $4R + 5N \leq 5610$ (frame manufacturing hours)
 $1.5R + 2.0N \leq 2200$ (wheel and deck manufacturing hours)
 $1.0R + 0.8N \leq 1200$ (QA and packaging hours)
 $R, N = \text{integer}$
 $R, N \geq 0$

	Razor	Navajo	
Profit Contribution (\$/unit)	150	160	
Units to Make	840	450	Total Profit (\$)
			198000 [1]

	Resource requirements		Required (hours)		Available (hours)
	Razor	Navajo			
Frame Manufacturing	4	5	5610 [2]	<=	5610
Wheels and Deck Assembly	1.5	2	2160	<=	2200
QA and Packaging	1	0.8	1200	<=	1200

[1]
 =SUMPRODUCT(C9:D9,C10:D10)
 [2]
 =SUMPRODUCT(\$C\$10:\$D\$10,C14:D14)

Excel fails on reproducibility

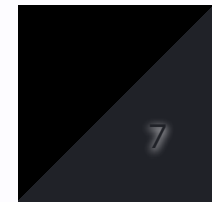
SDSS 2021, Soren Harner, RC2AI

6

Addressing Usability, Jim and team built

- *RC²* iOS and Mac Swift Client for Notebooks with remote computation, inspired by teaching at WVU
- *RSpark* a "Big Data" Docker-base compute environment with R, Spark, and other packages; used in short courses in Big Data with Spark

SDSS 2021, Soren Harner, RC2AI



Late nights discussing

- R vs. Python
- R Markdown vs. Jupyter
- Spark vs. Apache Arrow
- Tensorflow vs. PyTorch
- Stan vs. PyMC3
- Mac vs Linux
- Functional programming

- Docker, Kubernetes, GPUs

SDSS 2021, Soren Harner, RC2AI



JIM HARNER

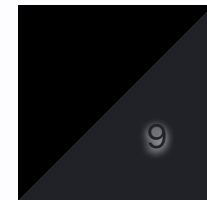
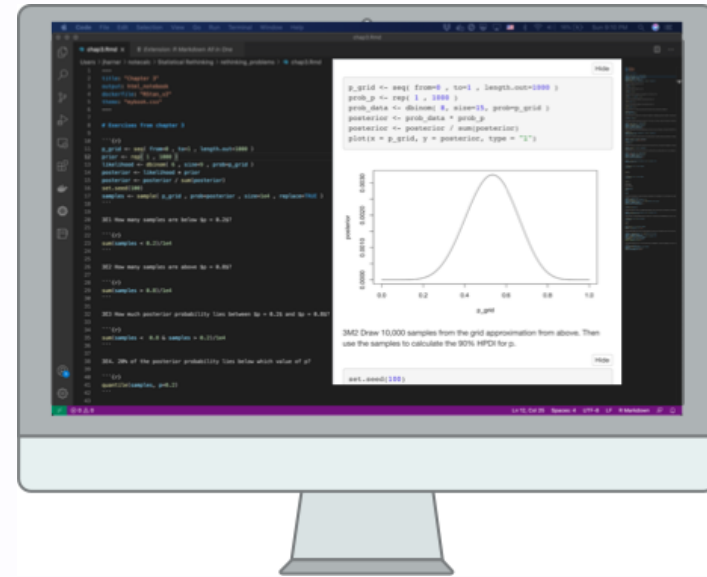


SOREN HARNER

Didact

- Author in markdown
- Add code chunks
- VS Code ecosystem
- Language servers
- Javascript interactivity
- CSS formating
- Pull request to publish

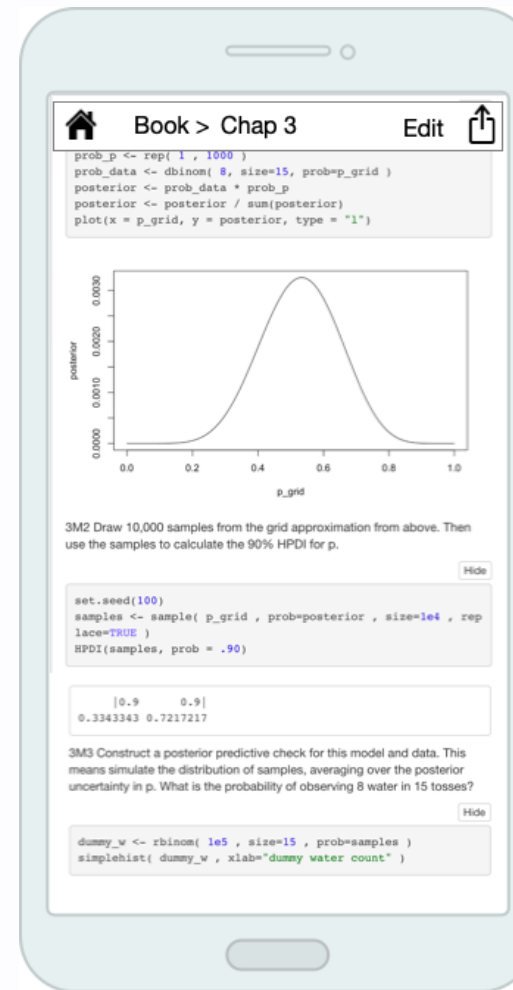
SDSS 2021, Soren Harner, RC2AI



NoteCalc

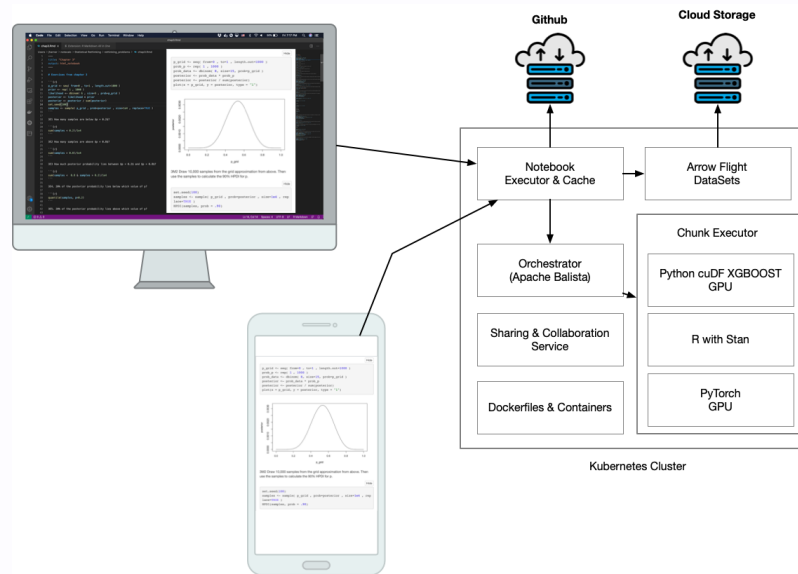
- Read, share, re-use interactive notebooks
- Mobile and Web
- Obviate the calculator
- Better than static books and reports
- Verifiable, tamperproof

SDSS 2021, Soren Harner, RC2AI



10

Remote computation and cloud services



SDSS 2021, Soren Harner, RC2AI

11

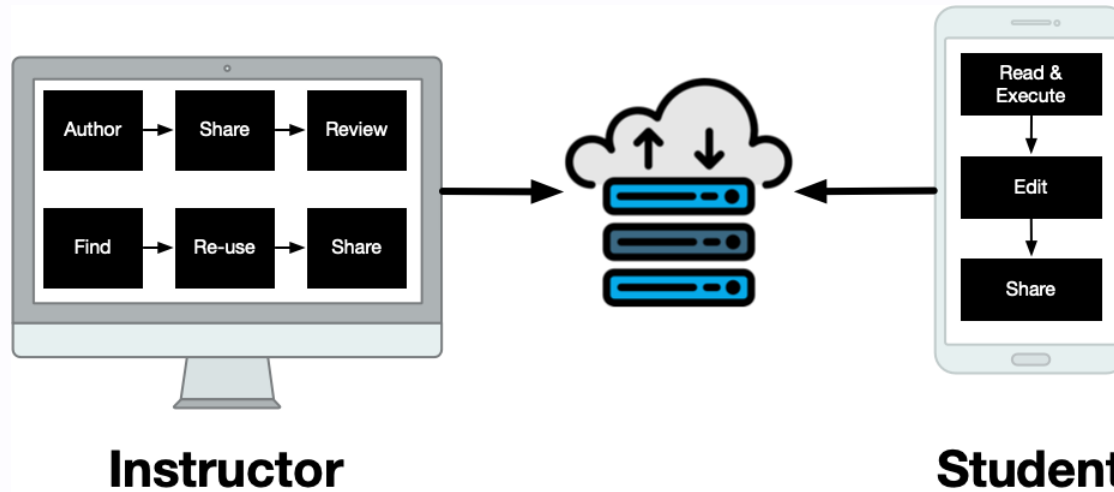
Reproducibility key to sharing

- Accessible code, e.g. Github
- Accessible, tamper-proof dataset vault, e.g. Arrow
- Accessible environment, e.g. Dockerfiles
- Notebook compiled into computational graph, i.e. shared random seed, externalized jobs
- Verifiable identity, auditable

SDSS 2021, Soren Harner, RC2AI



Instructor, student collaboration



SDSS 2021, Soren Harner, RC2AI

13

References

Introduction to Reproducible
Science in R

1,500 scientists lift the lid on reproducibility

What does research reproducibility mean?

Trust in science, social consensus and vaccine
confidence, Nature 2021 Survey

SDSS 2021, Soren Harner, RC2AI

14