# Lowering barriers to using computational notebooks

Tanaka Chitete
Bachelor of Science (Honours), Computer Science

Dr. Rahul Gopinath
Lecturer, School of Computer Science

THE UNIVERSITY OF SYDNEY

# Problems

# Setting up

- Problems are faced before computation has even begun
- Loading and cleaning data is menial

# Exploring and analysing data

- Exploration feels incomplete on many platforms
- Users find themselves stuck copying and pasting code

# Managing code

- Dependency management is difficult for users without software engineering expertise

- Such users are plunged into "dependency hell"

# Ensuring reliability

- Execution is unreliable and scaling to big data is unsupported
- Kernel crashes often result, causing data inconsistencies
- The lack of processing power is the main culprit

# Archiving

- Some notebooks are made with long-term use in mind
- Most platforms don't offer sufficient versioning or searching

# Securing

- Mitigating data leaks and controlling access are key concerns
- A lack of security measures forces custom implementations
- Such implementations are error-prone

# Sharing and collaborating

- Sharing happens *interactively*, *statically*, and as *read-only* files
- Collaboration via real-time editing is desired
- Existing sharing functionalities are insufficient
- Synced notebooks present coordination challenges

# Reproducing and reusing

- Marking work and furthering research hinge on reproducibility
- Users want to adapt and reuse code to save time
- Neither reproducing computation or code is feasible
- Notebook and machine dependencies are contributing factors

# Deploying to production

- Notebooks are better suited to "quick-and-dirty" work
- Cell-based programming doesn't coincide with best practices
- Deployment requires significant work

# Background

# Computational notebook

- Document that encapsulates code, text, and visualisations
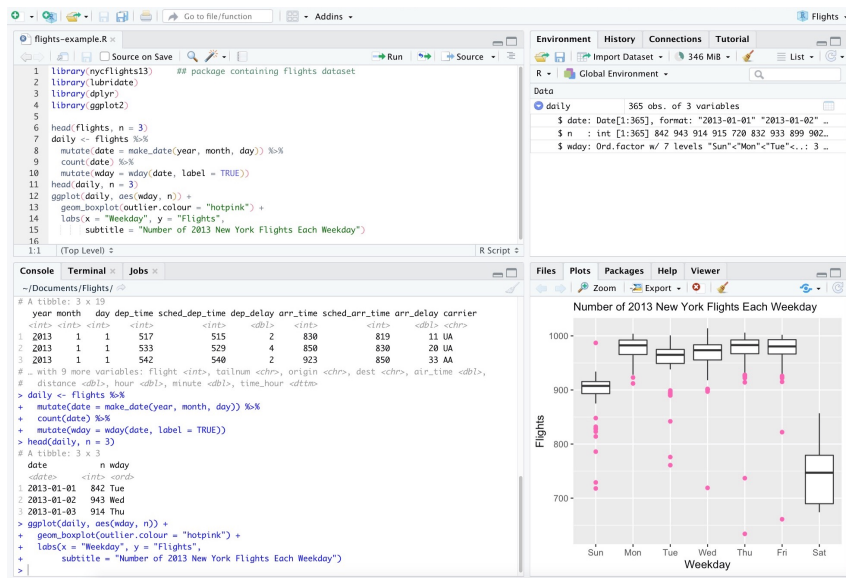- Organised as a sequence of cells

# Literate programming

- Paradigm at the foundation of the computational notebook
- Improves readability by considering programs as "works of literature"

# Computational platform

- Underlying system that runs the notebook
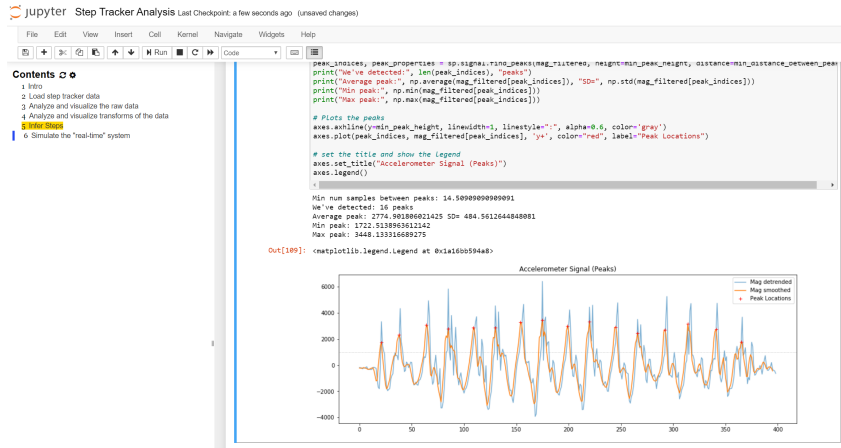- Provides the necessary compute to execute code
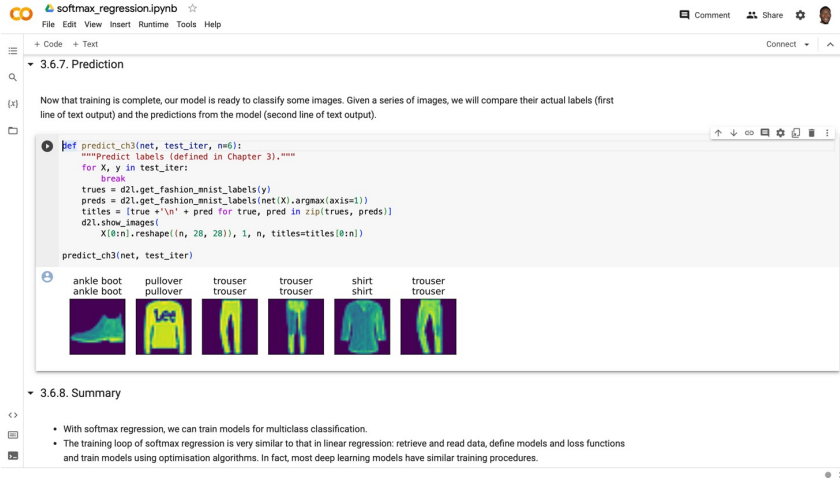
# Chronology

# RStudio



- Released in 2011 to offer a native platform for R

- Integrates editor, plotting, console, and variable views

- Features code-completion and syntax highlighting

- Integrates R library-set

- Provides Python interface

# Jupyter Notebook



- Unveiled in 2015 to offer a more versatile platform

- Distinguished by its kernel-based architecture

- Supports over 40 programming languages

- Considering Python alone, supports all major packages

# Colab



- Launched in 2017 to enhance ML education

- Based on Jupyter platform

- Offers free GPU access

- Supports Python, integrating major ML packages

- Connects to Google Drive

# Solution

# Demonstrating Yuna

`https://tanaka-chitete.github.io/yuna/`

# Implementing the interface

- CodeMirror supplies input cells
- Bootstrap supplies toolbar, buttons, and styling

# Executing code

- Pyodide powers code execution
  - Ports Python 3.11 runtime using Wasm
  - Loads runtime into the browser using JavaScript
  - Maintains standard library as a zip archive

# Loading packages

- Custom Python script embeds requirements
- micropip installs packages

# Restoring sessions

- cloudpickle serialises and deserialises all user-defined objects
  - Variables
  - Functions
  - Imports

# Parsing text

- Marked powers rich-text documentation

# Evaluation

# Managing code – Managing dependencies

- Installing, updating, and handling packages
- Platforms should offer integrated dependency management
- Yuna embeds all dependencies within the file, fully automating dependency management

# Ensuring reliability – Executing

- Running code, either select cells or the entire notebook

- Platforms must ensure execution is constantly-available

- Yuna embeds both the kernel and parser within each environment, guaranteeing availability
  - An environment can be as small as 27 MB

# Archiving – Versioning

- Tracking changes to source code files

- Notebooks rely on a wide range of artefacts

- Platforms must offer an integrated version control mechanism

- Yuna embeds the notebook, data, packages, and compute within each environment

# Sharing and collaborating – Sharing the notebook

- Making the notebook accessible to others
- If sharing interactively is infeasible, platforms should offer cell-collapsing, exporting, and comment-only viewing
- Yuna embeds the notebook and platform within each environment

# Sharing and collaborating – Sharing related artefacts

- Making the associated session, data and packages available
- Platforms should provide integrated artefact packaging
- Yuna embeds both data and packages within each environment

# Reproducing and reusing – Reproducing

- Reproducibility is the ability to consistently replicate results
- Platforms should be designed to inherently support reproducibility
- Yuna embeds all artefacts within each environment file

# Reproducing and reusing – Reusing and adapting

- Leveraging prior code for use in differing contexts
- Platforms should provide integrated dependency management
- By embedding dependencies within each environment, Yuna allows the user to simply copy the file

# Discussion

# Future work

- Loading data
- Usability testing
- Serialising
- Visualising
- Minimising file sizes
- Maintaining software
- Exporting contents

# Conclusion

# Contributions

- Compendium of domain knowledge
- 19-point chronology of the most significant platforms
- Yuna, a fully self-contained platform, soothing:
  - Managing code
  - Ensuring reliability
  - Archiving
  - Sharing and collaborating
  - Reproducing and reusing
- Technical evaluation of 13 state-of-the-art platforms

# Summary

## Demonstrating Yuna

`https://tanaka-chitete.github.io/yuna/`

## Reproducing and reusing – Reproducing

- Reproducibility is the ability to consistently replicate results
- Platforms should be designed to inherently support reproducibility
- Yuna embeds all artefacts within each environment file

## Future work

- Loading data
- Usability testing
- Serialising
- Visualising
- Minimising file sizes
- Maintaining software
- Exporting contents

## Contributions

- Compendium of domain knowledge
- 19-point chronology of the most significant platforms
- Yuna, a fully self-contained platform, soothing:
  - Managing code
  - Ensuring reliability
  - Archiving
  - Sharing and collaborating
  - Reproducing and reusing
- Technical evaluation of 13 state-of-the-art platforms

# References

Aashita. 2023. Beginners Guide to Databricks: Batch Processing and Streaming. *Databricks*.

cdhowe. 2021. The RStudio IDE is a free and open source integrated development environment (IDE) for R. It helps programmers develop software written in R, Python, and several other languages. *Wikimedia Commons*.

Souti Chattopadhyay, Ishita Prasad, Austin Z. Henley, Anita Sarma, and Titus Barik. 2020. What's wrong with computational notebooks. *In Proceedings: 2020 Conference on Human Factors in Computing Systems*.

Databricks. 2022. Databricks | Notebook Development Overview. *YouTube*.

Databricks. 2023a. Collaborate using Databricks notebooks. *Databricks*.

Databricks. 2023b. Databricks for Python developers. *Databricks*.

Databricks. 2023c. Databricks for R developers. *Databricks*.

Databricks. 2023d. Databricks for Scala developers. *Databricks*.

Databricks. 2023e. MLflow guide. *Databricks*.

Databricks. 2023f. SQL language reference. *Databricks*.

Databricks. 2023g. TensorBoard. *Databricks*.

Databricks. 2023h. What is Delta Lake? *Databricks*.

Databricks. n.d.a. Databricks is the data and AI company. *Databricks*.

Databricks. n.d.b. The best data platform is a lakehouse. *Databricks*.

Dave Gershgorn. 2017. Nerds rejoice: Google just released its internal tool to collaborate on AI. *Quartz*.

# References

Google. n.d. Welcome to Colab! *Google*.

Donald Knuth. 1984. Literate Programming. *The Computer Journal*.

Makeability Lab. n.d. Jupyter Notebook showing visualization of a 3-axis accelerometer to infer step counts. *Makeability Lab*.

Matplotlib. n.d. Image tutorial. *Matplotlib*.

Microsoft. 2023a. Train TensorFlow models at scale with Azure Machine Learning. *Microsoft*.

Microsoft. 2023b. Train PyTorch models at scale with Azure Machine Learning. Microsoft.

Microsoft. n.d. Azure Machine Learning. *Microsoft*.

Ron Miller. 2015. Microsoft Officially Launches Azure Machine Learning Platform. *TechCrunch*.

NumPy. n.d. NumPy: The Absolute Basics for Beginners. *NumPy*.

pandas. n.d. Intro to Data Structures. *pandas*.

Fernando Perez. 2014. Project Jupyter. *Speaker Deck*.

Project Jupyter. n.d. Jupyter: Free software, open standards, and web services for interactive computing across all programming languages. *Project Jupyter*.

RStudio. n.d.a. Python. *RStudio*.

RStudio. n.d.b. RStudio Open-Source Packages. *RStudio*.

RStudio Team. 2011. RStudio, a new open-source IDE for R. *Posit Software*.

Kevin Ushey. 2023. Code Completion in the RStudio IDE. *Posit Support*.