



Writing clean, responsible code



Why Clean Code?

- Switch mindset
 - coding by yourself -> coding collaboratively.
 - code living on your laptop -> code hosted live and automated
- Reproducibility
 - Someone trying to help you debug an error on Stackoverflow
 - You working on a project with multiple people
 - You receiving legacy code from someone else, or passing on your work to your successor
- Scalability
 - Just because it works on 1 csv file doesn't mean it won't run forever for the full dataset
 - Writing high performance code means less computing power, less memory issues, and less \$\$!



Technical Debt

- Technical Debt is a concept in software development that reflects the implied cost of additional rework caused by choosing an easy solution now instead of using a better approach that would take longer. (Wikipedia)
- Harsh deadlines can drive you to write sloppy, stream of consciousness code, but the hidden cost (technical debt) incurred of that can result in larger failures in the future when your code is automated.



Tenets of Clean Code

- Code is simple, clear, and concise.
 - Naming conventions should be explicit
 - Functions do one thing and one thing only. (e.g. “get_data”, “train_model”)
 - Do not repeat yourself! (If you find that you are repeating yourself, look to a loop)
 - Try to catch errors (e.g. `try`) instead of throwing default error codes
- Format correctly. Follow the style guides:
 - Pythonic code: <https://docs.python-guide.org/writing/style/>
 - PEP8 guide for Python: <https://www.python.org/dev/peps/pep-0008/>
- Code is clearly documented and commented.



Tenets of Clean Code

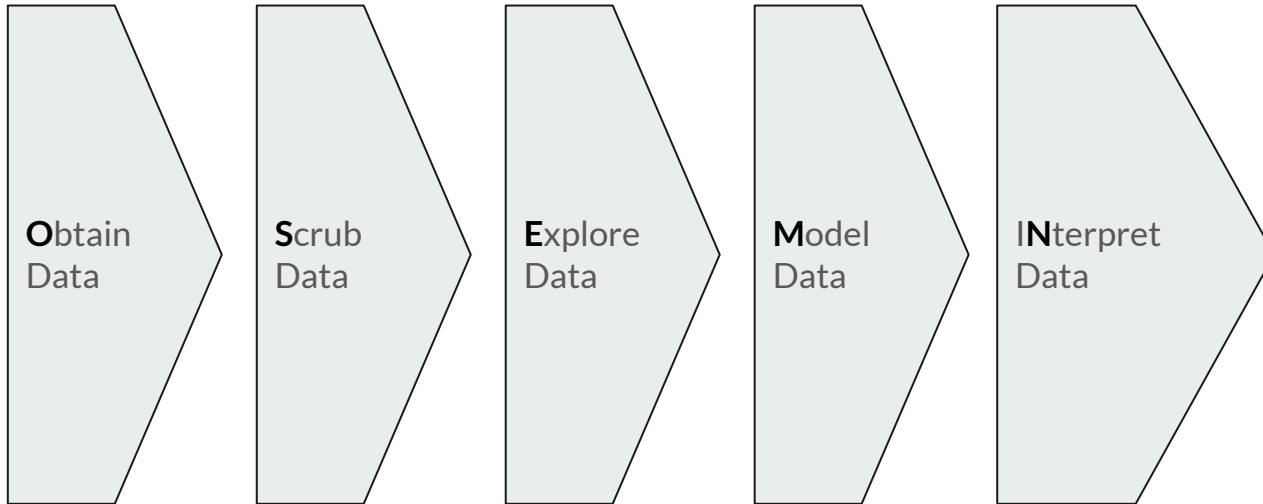
- The code is easily testable.
 - A testing unit (unit test) should focus on one tiny bit of functionality and prove it correct.
 - All tests are easy to understand and easy to change. Avoid “glamor tests”.
 - More here: <https://docs.python-guide.org/writing/tests/>
- Code review.
 - Have a friend, classmate, coworker review your code. Doesn't matter if that person has more, less, or equal experience than you.
 - Working alone? Try the rubber duck debugging method (https://en.wikipedia.org/wiki/Rubber_duck_debugging)
 - Learn version control! Git all the things!!
- Read the book [Clean Code](#) for more details and check out this [guide](#) specific for Python data work.



Complexities in your data workflow

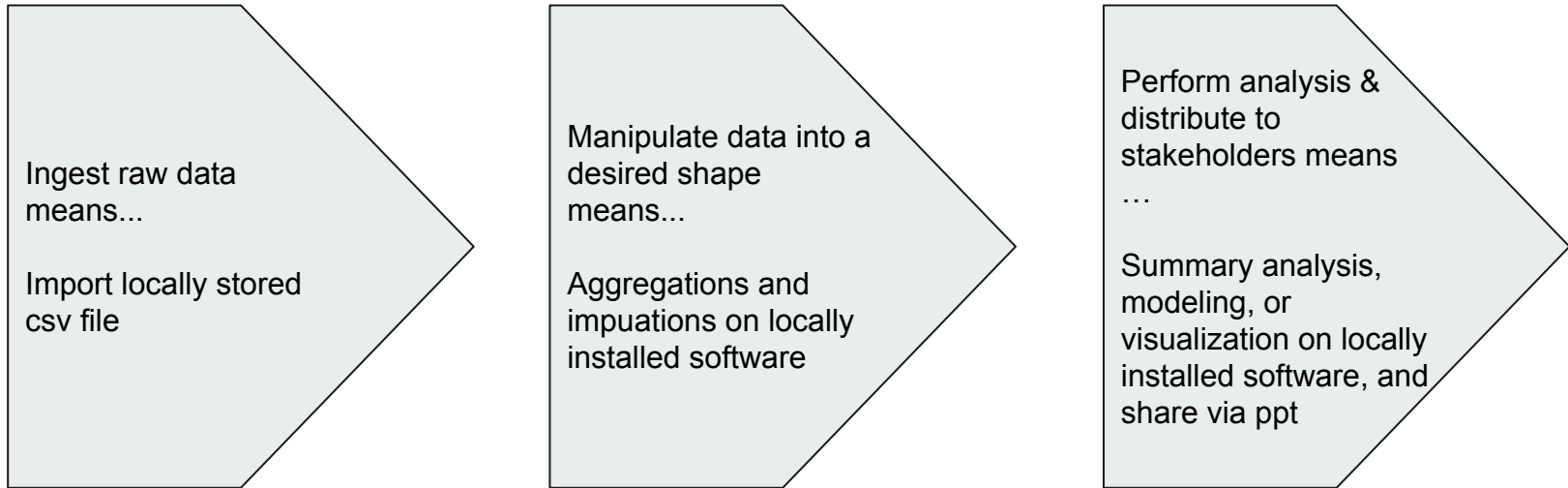


The OSEMN framework for data workflow





A simple example



But things can get complicated very quickly...



What if the data is more than a simple text file?

Ingest raw data means...

Use tools to parse data into **tabular form**.

Design a **data model**, choosing a local database client and perform **ETL**.

Manipulate data into a desired shape means

...

Query and join the data into a dataframe object fit for analysis.

Perform analysis & distribute to stakeholders means

...

Summary analysis, modeling, or visualization, and share via ppt



What if new data is periodically available?

Ingest raw data means...

The database needs to be on a **network** instead of local. **Jobs** needs to be written to ingest data every time it comes in.

Manipulate data into a desired shape means ...

Start building **data lakes** and **data marts** so that data is always ready for analysis.

Perform analysis & distribute to stakeholders means ...

Instead of static reports, consider online dashboards and deploying models.



What if every step is too large to be handled on a single machine?

Ingest raw data means...

Moving database to a **network**, and setting up infrastructure for **multi user access**.

Manipulate data into a desired shape means ...

Spin up **clusters** to utilize multiple machines. Optimize the best way to do every aggregation because efficiency really counts now!

Perform analysis & distribute to stakeholders ...

Spin up clusters to utilize multiple machines. Need to implement testing and make sure things **fail gracefully**.



What if stakeholder wants real-time data access and you want to save money?

Ingest raw data means...

Real-time stream processing using cloud computing services.

Manipulate data into a desired shape means ...

Real-time stream processing using cloud computing services.

Perform analysis & distribute to stakeholders means ...

Real-time model hosting using cloud computing services.



Major players in cloud computing ...

[Amazon Web
Services](#)

[Google Cloud
Platform](#)

[Microsoft Azure](#)

More on the design of data architecture? Check out the following two presentations:

[here](#) and [here](#)



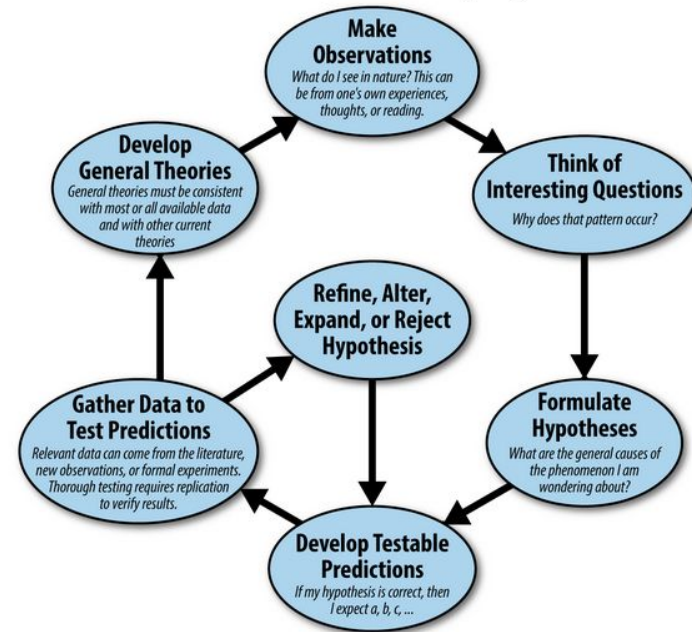
Complexities in working on a data team

A typical workflow for the solo data scientist ...

.. can be iterative because data scientists tend to come from an academic background and are used to the scientific method.

This works if it's solo work, but in a team, this workflow will change.

The Scientific Method as an Ongoing Process



A typical project team ...

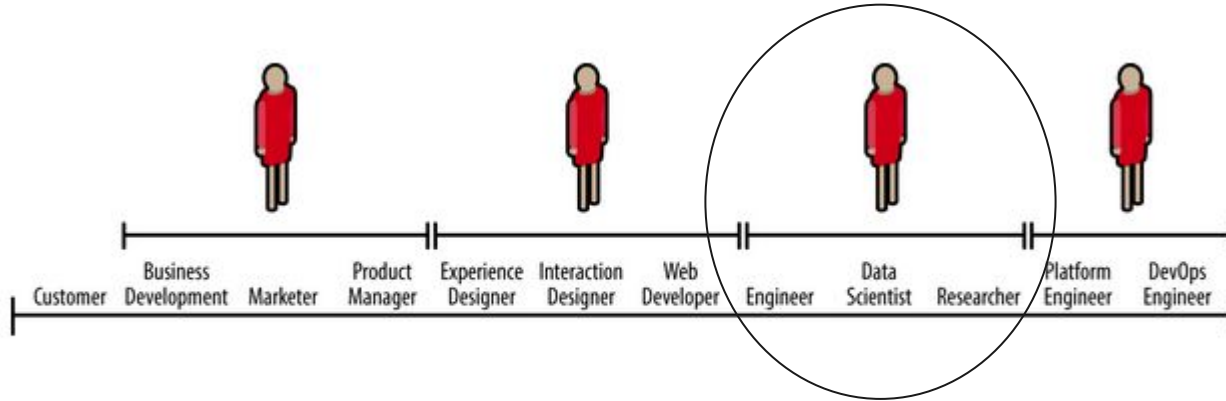
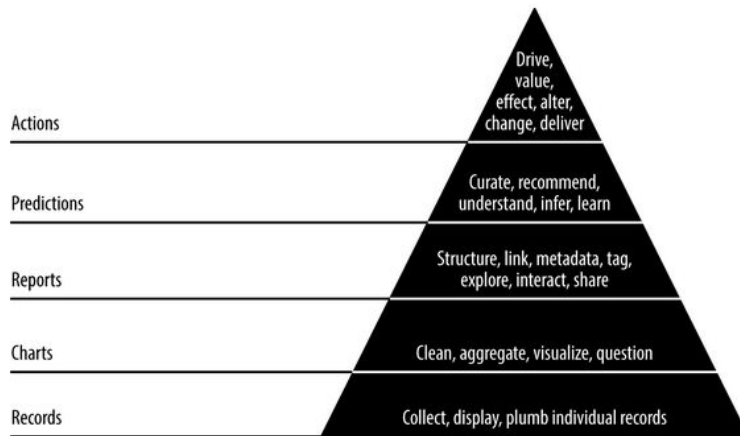


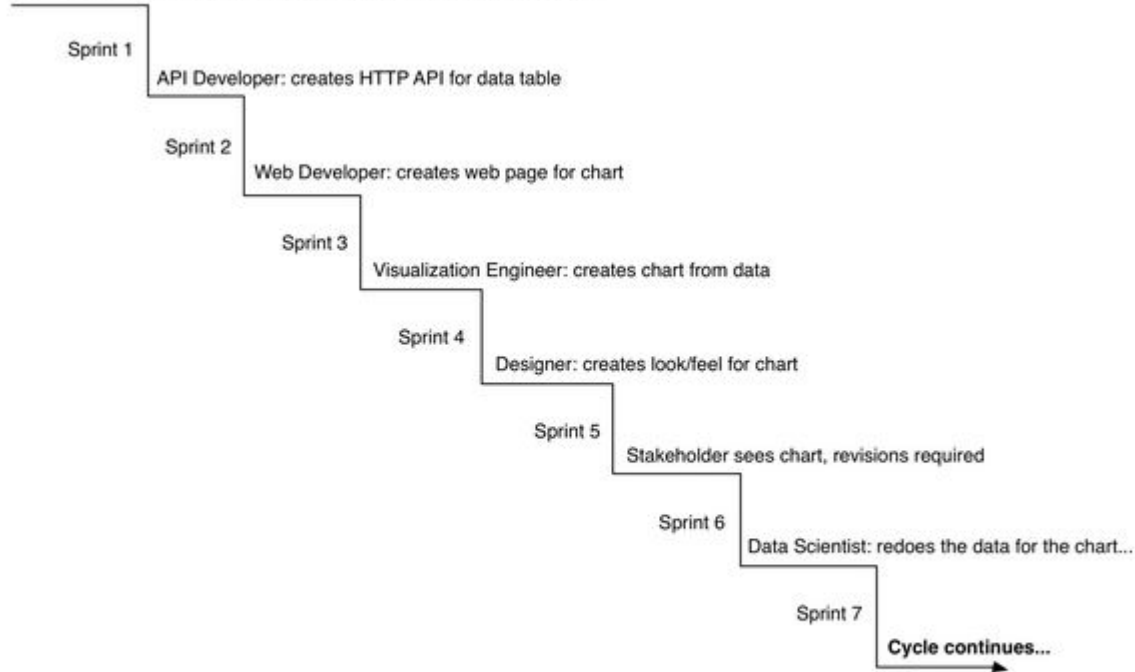
Figure 1-9. Broad roles in an Agile Data Science team

The Agile Data Science Manifesto

- Iterate, iterate, iterate: tables, charts, reports, predictions
- Ship intermediate output. Even failed experiments have output.
- Get meta. Describe the process, not just the end state.
- Climb up and down the data-value pyramid as we work.



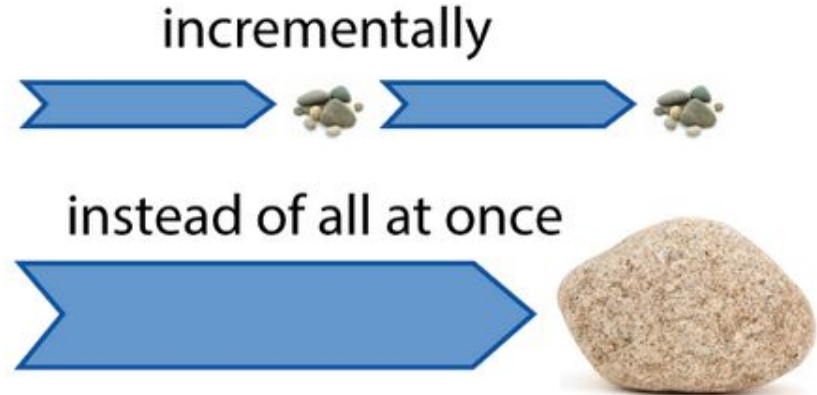
When you are a data scientist on a team ...



What is Agile?

Agile is a time boxed, iterative approach to software delivery that builds software incrementally from the start of the project, instead of trying to deliver it all at once near the end.

It works by breaking projects down into little bits of user functionality called user stories, prioritizing them, and then continuously delivering them in short 1-2 week cycles called iterations.





User Story

User stories are short, simple descriptions of a feature told from the perspective of the person who desires the new capability, usually a user or customer of the system. They typically follow a simple template:

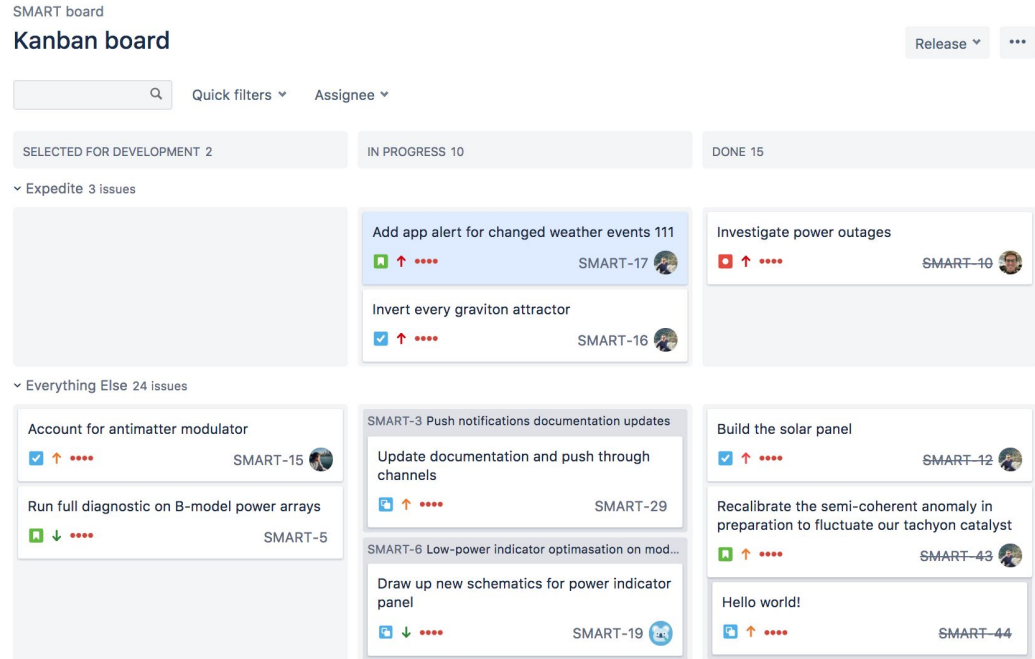
As a < type of user >, I want < some goal > so that < some reason >.

As a **data analyst**, I want **data in a certain format** so that **I can train a model on it**.

These user stories are then converted into bite-sized “tasks” with clear direction on how to accomplish this goal and then assigned to someone on the team.

Task Board

The **task board** is a visual display of the progress of the team during a **sprint**. It presents a snapshot of the current sprint backlog allowing everyone to see which tasks remain to be started, which are in progress and which are done.





Any more questions? Unmute yourself and ask!