



LEIDEN UNIVERSITY MEDICAL CENTER

Good Research Practice for computational data analysis

(part one)

Jeroen F. J. Laros

Leiden Genome Technology Center

Department of Human Genetics

Center for Human and Clinical Genetics



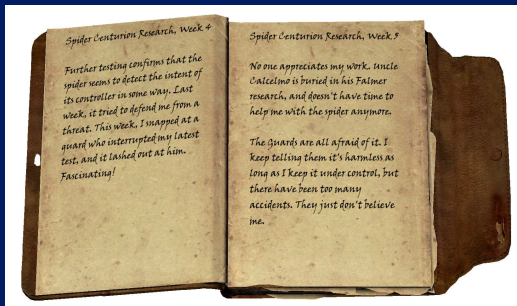
Good research practices

Figure 1: Lab journal.

In the lab: lab journals, standard operating procedures, etc.

Current practice

Pipeline development and data analysis have become increasingly important.

Surprisingly, there are no standard operating procedures for data analysis.

Current practice

Pipeline development and data analysis have become increasingly important.

Surprisingly, there are no standard operating procedures for data analysis.

Example.

- Person A creates a program.
- Person B modifies the program.

Current practice

Pipeline development and data analysis have become increasingly important.

Surprisingly, there are no standard operating procedures for data analysis.

Example.

- Person A creates a program.
- Person B modifies the program.
- Person B received PhD.
- Person A mails the program to person C.
- Person C can not reproduce the results of person B.

Current practice

Figure 2: Final resting place for data.

Even worse, the results disappear into a drawer.

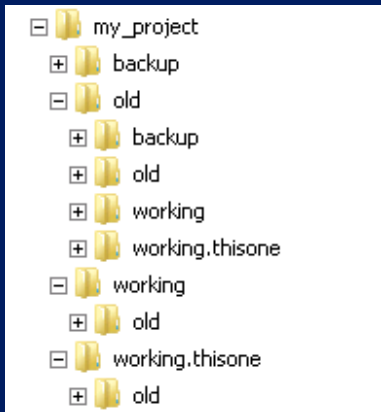
Current practice

Figure 3: “I have my own system.”

Current practice in our departement

Luckily, we have tackled these problems already.

Investments in infrastructure:

- Centralised storage.
 - Backup.
- Centralised computing.

Current practice in our departement

Luckily, we have tackled these problems already.

Investments in infrastructure:

- Centralised storage.
 - Backup.
- Centralised computing.

Adoption of *software engineering* solutions.

We want to share this knowledge with our colleagues in other departments.

“Beleidsinitiatief” (structural funding)

Proposal for a change in policy.

- The department Human Genetics has some nice solutions available.
- We always keep the big picture in mind.
- The LUMC as a whole can benefit.

Proposal was granted.

- One FTE structural.
- One additional FTE for one year.

M. Roos, P.A.C. 't Hoen

“Beleidsinitiatief” (structural funding)

The proposal consists of two main parts.

Work package 1.

- Version control.
- Interactive computational environment.

Work package 2.

- Data stewardship.

“Beleidsinitiatief” (structural funding)

The proposal consists of two main parts.

Work package 1.

- Version control.
- Interactive computational environment.

Work package 2.

- Data stewardship.

Work package 1 and 2.

- Education.

Git

The management of changes to documents, computer programs, large web sites, and other collections of information.
— Wikipedia.

<http://www.git-scm.com/>

<https://github.com/>

Git

The management of changes to documents, computer programs, large web sites, and other collections of information.
— Wikipedia.

General features:

- Keeping track of your files in an orderly manner.
 - Hiding old versions.
 - Recording who made changes and when.
- Enables collaboration.

<http://www.git-scm.com/>

<https://github.com/>

Why version control?

For a single user:

- Revert files to a previous state.
- Revert the entire project back to a previous state.
- Review changes made over time.
- Backup.

Why version control?

For a single user:

- Revert files to a previous state.
- Revert the entire project back to a previous state.
- Review changes made over time.
- Backup.

For multiple users:

- A reliable way to share files between people/computers.
- Allow multiple people working on the same project at the same time.
- Conflict resolution.
- See who made which changes at which time.

Collaboration

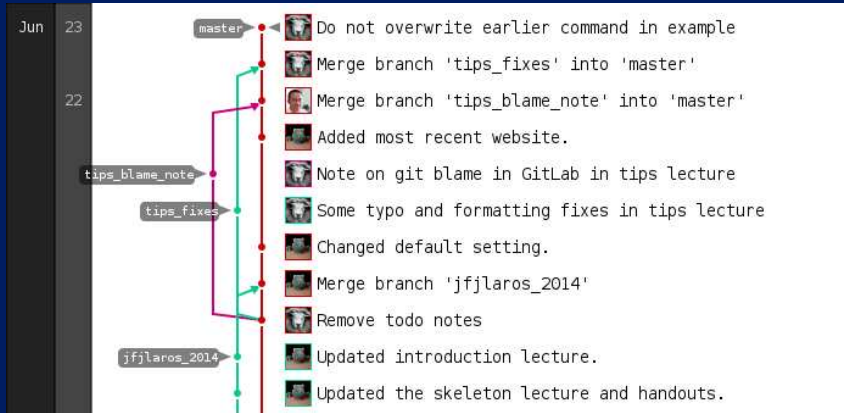


Figure 4: Collaboration with many people.

Tracking of changes

257	295		self._set_field(member, 1)
258	296		self._set_field(member, 1, 'FLAGS_2', _bit)
259	-	-	self._set_field(member, 26)
	297	+	self._set_field(member, 4)
	298	+	
	299	+	self._parse_crossover(member['ID'])
	300	+	
260	301		self._set_field(member, 1, 'FLAGS_3', _bit)
261	302		self._set_field(member, 205)
262	303		

Figure 5: Compare versions.

Not limited to the previous version and the latest one.

- Different authors.
- Any two versions.

Documentation

master fam-parser / fam_parser / +

Name	Last Update	Last Commit	History
..			
README.md	14 days ago	Jeroen F.J. Laros	Updated documentation.
container.py	10 days ago	Jeroen F.J. Laros	Documentation.
fam_parser.py	10 days ago	Jeroen F.J. Laros	Added support for the + and - flags.

README.md

Development of the FAM parser

Preparations

First, make sure that `wine` is installed and that your CPU allows execution of 16-bit instructions:

```
echo 1 > /proc/sys/abi/ldt16
```

Figure 6: Documentation and programs in one place.

Projects

Most of us work on multiple projects with multiple people.

That is why it is convenient to:

- Have everything in one place.
 - Data.
 - Code.
 - Documentation.
- Have the same structure for all projects.

Projects

Most of us work on multiple projects with multiple people.

That is why it is convenient to:

- Have everything in one place.
 - Data.
 - Code.
 - Documentation.
- Have the same structure for all projects.

This also makes transferring projects easier.

Projects

Ideally, every directory in the project has a README file.

directory	description
data	Raw immutable data.
doc	Sample sheets, papers, etc.
src	Programs specific for this project.
analysis	Results.

Table 1: Project layout.

This structure is used by the LGTC, later adopted by SASC.

<https://git.lumc.nl/lgtc-bioinformatics/project-skeleton/>

Interactive computational environments

Combine code execution, text, mathematics, plots and rich media into a single document.

Ideal for exploration of data.

- Documentation and code are interwoven.
- Results are displayed inline.
- Web based.
- Versions.

<http://ipython.org/notebook.html>

Interactive computational environments

Combine code execution, text, mathematics, plots and rich media into a single document.

Ideal for exploration of data.

- Documentation and code are interwoven.
- Results are displayed inline.
- Web based.
- Versions.

Integration with GitLab.

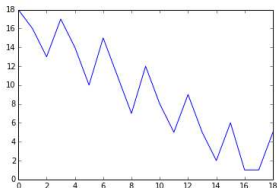
<http://ipython.org/notebook.html>

iPython notebook

Ranked data

```
In [7]: def ranked(data):
        rank = {v: k for k, v in enumerate(sorted(data))}
        return map(lambda x: rank[x], data)

        plot(ranked(data));
```



By using ranks instead of values, we circumvent a lot of scaling issues. Notice that this plot looks a lot like the logarithmic plot above, except that the "dampening effect" is gone.

Figure 7: *iPython notebook*.

Overview

We offer the following courses:

- Introduction to clusters.
- Introduction to Git.
- Python programming.

Apart from this, we:

- Teach people running our pipelines (SASC).
- Answer questions about the infrastructure.
- Keep documentation up to date.

<https://humgenprojects.lumc.nl/>

Introduction course clusters

Half day course.

An overview of the available infrastructure.

In particular:

- Connecting to a cluster.
- Using the Sun Grid Engine.
- Do's and don'ts.
- Makefiles.

Git

Everyone in the Bioinformatics field:

- Software development.
- Project management.
- Collaboration.

M. Vermaat, W. Arindrarto

<https://humgenprojects.lumc.nl/trac/humgenprojects/wiki/git>

Git

Everyone in the Bioinformatics field:

- Software development.
- Project management.
- Collaboration.

Topics:

- Git Basics
- Branching
- Remotes
- Project skeleton / git annex

M. Vermaat, W. Arindrarto

<https://humgenprojects.lumc.nl/trac/humgenprojects/wiki/git>

Programming in Python

Four day course.

- Python basics.
- Standard data structures.
- Working with NumPy arrays.
- Plotting with matplotlib.
- Object-oriented programming.
- The Biopython library.

M. Vermaat, W. Arindrato, Z. Tatum, W.Y. Leung

<https://humgenprojects.lumc.nl/trac/programming-course>

Data stewardship

This topic will be covered by Marco Roos at a later time.



Acknowledgements:

Marco Roos
Peter-Bram 't Hoen
Martijn Vermaat
Zuotian Tatum
Wibowo Arindrarto
Wai Yi Leung
Michel Villerius
Madeleine Nivard
Silvère van der Maarel