# Software development bestpractises for reproducible research

Alastair Porter

20 September 2018

MIP-Frontiers Meeting

### I'm Alastair

- \* I'm a researcher in the MTG
- I do research, but the majority of my time is spent writing software
- I have experience working in large organisations, and in community driven projects with hundreds of thousands of users
- This talk is based on my experience writing software in a research context, and questions that I get asked by students

# The ideal research pipeline

Researcher A ("Producer")

- Read background papers
- Do own research
- Publish paper X

Researcher B ("Consumer-Producer")

- Read paper X
- Understand/reproduce results in paper X
- Do more research building on X
- Publish paper Y that cites X / produce product that uses X

Plumbley et al. 2012. Research reproducibility workshop

# The real research pipeline

Researcher A ("Producer")

- Read background papers
- Do own research (including lots of coding)
- Publish paper X (no source code or data. May describe the algorithm incorrectly)

Researcher B ("Consumer-Producer")

- Read paper X
- Can't reproduce or use results in paper X
- Tear out hair
- Give up / do something else

Plumbley et al. 2012. Research reproducibility workshop

### Research is software development

Papers need code

- \* How many of you are programming as part of your research?
- \* How many of you did an undergraduate degree in CS or Software Development?
- \* How many hadn't programmed before you started here?

### Overview of this talk

- I. Code layout and organisation
- 2. Testing
- 3. Documentation
- 4. Source control
- 5. Data
- 6. Licensing
- Not a tutorial talk.
  - \* You will have to take these ideas and investigate them further
- \* Time for questions at the end, but please stop me if you want to discuss something

### A small disclaimer

# I'm not very good at this

- 2014 & 2015: I gave a talk similar to this one, about best practises in making research reproducible
- \* 2015: Did an experiment submitted a paper to a conference
- Paper was rejected
- 2016: I wanted to re-visit the paper, change some processes, and re-submit

 Should be easy, right? Just read the documentation and run the right script

### No

I didn't even know what file to run to reproduce the results, let alone where to start making changes

### Listening to a talk will not make you better

- \* That is to say, this is not magic
- You can't just decide to write good software and suddenly it will be good
- Like everything else, pick a small change to make, practise it until it's second nature. Repeat

## There is a happy ending

### AudioCommons Jamendo tools

This project contains tools used to get metadata and recordings from Jamendo for use in the AudioCommons project.

We are currently working with the Jamendo Licensing catalog. This catalog has about 200,000 available tracks, of which we are using a subset of about 80,000.

#### Setup

This code uses python 3. Make sure to use a python3 virtual environment (virtualenv -p python3.4 ve) or use python3 / pip3.

Install dependencies with

\$ pip -r requirements

To configure, copy config.py.dist to config.py and fill in your jamendo api key.

#### Usage

To get the contents of the licensing catalog

\$ python getlicensing.py <outputname.json>

To update the catalog with entries that have been published since a playlist was downloaded

\$ python getlicensing.py -u <existingplaylist.json> <outputname.json>

To get the contents of a playlist

### I'm here to transfer knowledge

- Observations and things that I've learned from 10 years of programming and software development
- I will list some main points which have helped me through the years, why they're a good idea, and why they've worked for me
- Almost all of these points will require more up-front work.

# "But I'm a very busy person"

- Your first thought will probably be "But this is going to take much longer. I don't have time for this, I have a paper due really soon!"
- \* Additional effort up front will save you time in the future

### I'm a very busy person

HOW LONG CAN YOU WORK ON MAKING A ROUTINE TASK MORE EFFICIENT BEFORE YOU'RE SPENDING MORE TIME THAN YOU SAVE? (ACROSS FIVE YEARS)									
	HOW OFTEN YOU DO THE TASK								
		50/DAY	5/DAY	Daily	WEEKLY	MONTHLY	YEARLY		
	1 SECOND	1 DAY	2 HOURS	30 MINUTES	4 MINUTES	1 MINUTE	5 SECONDS		
	5 SECONDS	5 DAYS	12 HOURS	2 HOURS	21 MINUTES	5 MINUTES	25 SECONDS		
	30 SECONDS	4 WEEKS	3 DAYS	12 HOURS	2 HOURS	30 MINUTES	2 MINUTES		
HOU MUCH	1 1 MINUTE	8 WEEKS	6 DAYS	1 DAY	4 HOURS	1 HOUR	5 MINUTES		
TIME	5 MINUTES	9 MONTHS	4 WEEKS	6 DAYS	21 HOURS	5 HOURS	25 MINUTES		
SHAV	F 30 MINUTES		6 MONTHS	5 WEEKS	5 DAYS	1 DAY	2 HOURS		
	1 HOUR		10 months	2 MONTHS	10 DAYS	2 DAYS	5 HOURS		
	6 HOURS				2 MONTHS	2 WEEKS	1 DAY		
	1 DAY					8 WEEKS	5 DAYS		

https://xkcd.com/1205/

### Starting a new project

### \* You have an idea for a project. Great, let's start



https://www.flickr.com/photos/jeanbaptisteparis/724619122/

### Advice from the Internet

Follow



Micah Allen @micahgallen

My universal directory structure, forged from years of bad organization:

/Projects .../inProgress ...../ProjectName ...../docs ...../code ...../data ..../figures .../published .../submitted

### You're welcome.

9:20 AM - 29 May 2018



Q 47 t͡] 412 ♡ 1.7K ⊠



Replying to @micahgallen

Excellent. I use something similar with the important detail that folders are named with the date, year first: YYYY\_MM\_ProjectName. After a many years, remember every project name can be tricky but remembering about the time it happened is easier.

Follow

```
1:27 PM - 29 May 2018 from Florida, USA
```



https://twitter.com/micahgallen/status/1001362580710088704 https://twitter.com/ProfZare/status/1001424807429275649

# Code layout and organisation

80	SamaEstimation									
<	215 GB Volume UPFWor	k PhD	githubCode	SamaEstimation					٩	= :::
mode	S	×	SamaEstima	tion	×	pdf				×
Name							<ul> <li>Size</li> </ul>	Туре	Mod	ified
	getiOlstats.m						2.5 kB	Text	Mar	28 2013
	getMelodyFeaturesForSama_ch	unks.m					4.3 kB	Text	Mar	26 2013
	HarmonicSubtraction_New.m						7.1 kB	Text	Oct	17 2013
	hpEnhance.m						1.5 kB	Text	Feb	27 2013
	hpEnhance_chunks.m						4.6 kB	Text	Маг	25 2013
	hpEnhance_Nochunks.m						1.5 kB	Text	Mar	21 2013
	HPS_BatchProcess.m						2.2 kB	Text	Aug	1 2013
	HPS_BatchProcess_New.m						2.8 kB	Text	Aug	1 2013
	HPSS_PrelimExpt.m						1.8 kB	Text	Feb	27 2013
	HPSS_PrelimExpt_Batch.m						1.7 kB	Text	Mar	19 2013
	HPS_Thoshkahna.m						297 bytes	Text	Mar	21 2013
	HPS_Thoshkahna_chunks_Perc.r	n					710 bytes	Text	Mar	22 2013
	Jagadandakaraka_Expt.sv						50.7 kB	Text	Feb	27 2013
	LSTMNeuralNetworkExpts.m						5.0 kB	Text	Арг	11 2013
	lstmNNexpt_OneFold.m						2.3 kB	Text	Арг	11 2013
e	pitchExtract.py						508 bytes	Text	Aug	1 2013
0	pitchExtractBatch.py						1.1 kB	Text	Aug	1 2013
	PlotFeatures_All.m						12.1 kB	Text	Apr	3 2013
	Readme.txt						11 bytes	Text	Feb	27 2013

### Code layout and organisation

- \* what is hpEnhance\_chunks? \_nochunks?
- What is normal hpEnhance?
- \* What is BatchProcess\_New? What does it do new which is different?
- \* Don't copy an entire file and rename it to try something new
  - \* You know what these files do now, but will you in 6 months?
  - \* What if you find an error in your evaluation component? Will you remember to change the original code and test that too?

## Breaking up code into units

- \* Don't put all your code into one method in one file
- \* Why not? It makes it more difficult to reason about parts
- \* makes testing harder (more on this later)
- \* makes reuse harder

# Breaking up code into units

- Split code into smaller blocks
  - Load
  - \* Preprocess
  - Main algorithm
  - Calculate statistics
- Change the algorithm? Make a second function instead of copying all additional code too
- Have a runtime flag (or setting variable) which lets you choose which algorithm to use

# Breaking up code into units

\* How does this save time?

- \* More up-front work to split code into modules (although, not much more if you split it from the beginning)
- \* You only have one version of common code
- Save time having to look through each of your versions to find the most up-to-date one
- Common interface to all of your versions of code—Choose between them quickly for evaluation

## Hard-coded paths

20	<pre>def query():</pre>
21	g = rdflib.Graph()
22	<pre>mo = g.parse("C:\Users\Sergio\Dropbox\QMUL\Data\searchOntology.owl")</pre>
23	for subj, pred, obj in mo:
24	<pre>moSub.append(subj)</pre>
25	moPred.append(pred)
26	moObj.append(obj)
27	print "Subject: " + subj
28	print "Predicate: " + pred
29	print "Object: " + obj

- Not portable off your computer (very bad if someone else wants to run it)
- Use arguments instead

# Using arguments



### Pass arguments when you run your tool

Run/Debug Configurations						
Name: cool_project		<u>S</u> hare	Single instance only			
Configuration Logs						
Script:	/Users/alastair/code/cool_project/process.py					
Script parameters:	param_here		•			
▼ Environment						

\* Using an IDE? You can still add arguments

## Using arguments



### In python, don't read sys.argv

if \_\_name\_\_ == "\_\_main\_\_":
 parser = argparse.ArgumentParser(description="Process a datafile into parts")
 parser.add\_argument("-d", action="store\_true", help="Print debugging info")
 parser.add\_argument("input", help="Input json file")
 parser.add\_argument("result\_dir", help="Result directory to write items")
 args = parser.parse\_args()
 myfunction(args.input, args.result\_dir, args.verbose)

Use argparse

(but it's so complex...)

## Argparse magic

### \* But we want our work to be understandable

```
1.alastair@MacBook-Air:~/repro-research-2017$ python convert_metadata.py --help
usage: convert_metadata.py [-h] [-d] input result_dir
Process a datafile into parts
positional arguments:
    input Input json file
    result_dir Result directory to write items
optional arguments:
    -h, --help show this help message and exit
    -d Print debugging info
    alastair@MacBook-Air:~/repro-research-2017$
```

\* You get documentation for free

## Other options for Arguments

- Fire (<u>https://github.com/google/python-fire</u>)
- Magically creates help text and argument handling from a python object or model

Here's an example of calling Fire on a class.

```
import fire
class Calculator(object):
    """A simple calculator class."""
    def double(self, number):
        return 2 * number

if __name__ == '__main__':
    fire.Fire(Calculator)
```

Then, from the command line, you can run:

```
python calculator.py double 10 # 20
python calculator.py double --number=15 # 30
```



\* How does this save time?

- \* Less work when you want to release it publicly
- If you need to run the software on another machine, less work to configure it (pass all configuration options as arguments)
- \* Faster to deploy on the HPC

## Testing and Automation

- We have broken our code into small blocks
- \* How do we make sure the blocks work as they should?
- Tests!
- Are tests worth the effort?
- That depends on who you ask and what your definition of testing is

# "Writing tests takes too long"

- \* Yes, it does.
- In software development, it's not unusual that writing tests takes as long as writing the software. This means that your development time could **double**
- So, why is it good?

We are saving time when making changes or checking code

# Typical development workflow

- I. Make a change
- 2. Restart your python terminal
- 3. import your module
- 4. open a file
- 5. forget to import json module
- 6. import json. open file again
- 7. get a field from json file
- 8. call your method with data as an argument
- 9. print out result
- 10. see that the formatting of your output isn't right
- II. go to step I

## Easier development

- I. Make a change
- 2. run python mycode/test\_parser.py
- 3. Fix any errors
- 4. Repeat

## Testing as automation

Rule of thumb: If you're doing something manually more than
 2 or 3 times, perhaps you should write a test or automate it

- Yes, you could press the up button in ipython. It's still slower than writing a test
- If you come back to your code in 3 weeks time can you remember how you tested it?

## Testing as documentation

- A test is a kind of documentation for how to use a system (or a method)
- Sometimes when you write that documentation you realise the interface to your code is going to be unnecessarily clumsy
- Change your code to make it easier to automate and to document



\* How does this save time?

- Test only a small part of your algorithm instead of waiting for the entire load/process stage to complete
- Influences the design of small connectable blocks of code, potentially saving refactoring time in the future

## Testing example

\* There's a bug in this code

```
In [ ]: def filelist for makam and form(fileList, makam, form):
            '''Selects/returns a file list (a subset of the whole list: 'fileList')
            given the filters: makam(str) and form(str)'''
            selectedFiles = []
            for filepath in fileList:
                filename = os.path.basename(filepath)
                if filename.split('--')[0] == makam and filename.split('--')[1] == filename:
                    selectedFiles.append(filepath)
            return selectedFiles
        # Let's pick two makams only differing in seyir and plot melodic curves
        # for first sections in the sazsemaisi form
        makams = ['rast', 'mahur']
        form = 'sazsemaisi'
        # Example
        txtFilesList = ['muhayyer--bozlak--serbest--havayi da--asaf guven.txt',
                        'ussak--sarki--nimsofyan--gordum seni--sadettin kaynak.txt',
                        'mahur--sazsemaisi--aksaksemai----nikolaki.txt',
                         'sipihr--pesrev--muhammes----dilhayat kalfa.txt',
                         'mahur--sazsemaisi--aksaksemai--bahar 1--goksel baktagir.txt',
                        'buselik--sarki--duyek--sevgim bitti mi--arif sami toker.txt']
        for index, makam in enumerate(makams):
            selectedFileList = filelist for makam and form(txtFilesList, makam, form)
            print('Files in makam {} and form {}'.format(makam, form))
            print(selectedFileList)
```

## Testing example

- \* Make a separate test file with known inputs and outputs
- \* You don't need to wait for the rest of the cell to run

```
In [ ]: def filelist for makam and form(fileList, makam, form):
            '''Selects/returns a file list (a subset of the whole list: 'fileList')
            given the filters: makam(str) and form(str)'''
            selectedFiles = []
            for filepath in fileList:
                filename = os.path.basename(filepath)
                if filename.split('--')[0] == makam and filename.split('--')[1] == filename:
                     selectedFiles.append(filepath)
            return selectedFiles
        # Let's pick two makams only differing in seyir and plot melodic curves
        # for first sections in the sazsemaisi form
        makams = ['rast', 'mahur']
        form = 'sazsemaisi'
        # Example
        txtFilesList = ['muhayyer--bozlak--serbest--havayi da--asaf guven.txt',
                         'ussak--sarki--nimsofyan--gordum seni--sadettin kaynak.txt',
                         'mahur--sazsemaisi--aksaksemai----nikolaki.txt',
                         'sipihr--pesrev--muhammes----dilhayat kalfa.txt',
                         'mahur--sazsemaisi--aksaksemai--bahar 1--goksel baktagir.txt',
                         'buselik--sarki--duyek--sevgim bitti mi--arif sami toker.txt']
        for index, makam in enumerate(makams):
            selectedFileList = filelist for makam and form(txtFilesList, makam, form)
            print('Files in makam {} and form {}'.format(makam, form))
            print(selectedFileList)
```

### Documentation

- \* Great, you have some code. How does it work?
- \* Add a description of what the software does. What is the expected input format? Is there an example of what a successful run looks like?
- Is this the implementation of code which you published?
   After the paper is accepted, cite it!

### File documentation

### \_\_\_\_\_

This script processes the datafiles generated by the foo-project download script

The input should be a json file, formatted like this:

{"itemid1": {"field1": "value1", "field2": False}, "itemid2": {"field1": "value6", "field2": True}}

For each item key, a new file, `itemid.json` will be created in the output directory containing the contents of the value of this key.

### Project documentation

#### https://github.com/sertansenturk/ahenkidentifier

#### E README.md

build passing codecov 97% code climate 4.0 version 1.5.2 DOI 10.5281/zenodo.259699 License AGPL v3

### <sup>®</sup> ahenkidentifier

Identifies the ahenk (transposition) of a makam music recording given the tonic frequency and the symbol (or the makam)

If you are using ahenkidentifier in your work, please cite:

Şentürk, S. (2016). Computational Analysis of Audio Recordings and Music Scores for the Description and Discovery of Ottoman-Turkish Makam Music. PhD thesis, Universitat Pompeu Fabra, Barcelona, Spain.

#### Usage

from ahenkidentifier.ahenkidentifier import AhenkIdentifier

```
ahenk = AhenkIdentifier.identify(tonic_freq, makam)
# or
ahenk = AhenkIdentifier.identify(tonic_freq, tonic_symbol)
```

#### The inputs are:

# tonic\_freq : The frequency of the tonic in Hz.
# makam/tonic\_symbol: The algorithm can either accept the makam-slug or
# the note symbol of the tonic in SymbTr format as a string (e.g. B4b1).

For the makam-slug names, check the json file in the data folder. The slugs are the same with the ones in the filenames of the scores in the SymbTr collection. The tonic symbols are notated as [Note pitch-class][Octave] (Accidental Symbol(Holderian Comma)), e.g. B4b1

Why document?

\* How does this save time?

- Because you are going to forget what your code does
- Save yourself time having to read all of your files to work out how to run the software
- If you want to publish the software, you now have documentation without any extra work

### Source control



http://soundsoftware.ac.uk/why-version-control

### Do you remember what you did last semester?

- Your code works
- You make a change
- \* You make another change
- \* ...
- Copy your code to your laptop
- Make another change
- \* Your code doesn't work
- \* Easy way of keeping distributed backups (and knowing what is most up to date)

### Source control

Use git



Available for mac, windows, linux

- \* GitHub will give you a free place to store your code
- If you want to use mercurial or subversion, that's absolutely fine

### Academic Discount

https://education.github.com



STUDENT DEVELOPER PACK



### Get the Student Developer Pack

Dozens of free resources from great companies to help students learn.

Get the pack

### Tutorials



Almost all text on GitHub is processed through a markup language called *Markdown* — it's an easy way to include simple formatting (like *italics*, **bold words**, lists, and links).

GitHub Flow is a lightweight, branch-based workflow that supports teams and projects where deployments are made regularly. This guide explains how and why GitHub Flow

-

# Issues (Organise your life)

MTG / dunya	Add Repo	O Unwatch	- 16	★ Star	13	¥ Fork	7
<> Code () Issues 42 () Pull requests 0 () Projects 0	🗉 Wiki 🥠 Pulse	III Graph	s 🔅 S	Settings			
Filters - Q is:issue is:open	Milestones					New iss	ue
① 42 Open ✓ 367 Closed	Author -	Labels -	Mileston	es +	Assignee	- Sor	t-
Wrong tonic identification result returned in joint and #413 opened 28 days ago by sertansenturk Makam	lysis <mark>makam</mark>					Ļ	5
<ul> <li>Storing MBID redirects</li> <li>#410 opened on Nov 24, 2016 by sertansenturk</li> </ul>							
Missing audio features makam #408 opened on Nov 2, 2016 by sertansenturk Makam							
Missing external identifier for the recording Düşse Z #406 opened on Nov 2, 2016 by sertansenturk Makam	ülfünden (Hicaz İkino	ci Beste) mai	kam		3	Ļ	] 1
visualise mel scale in Dunya interface for lyrics-to-au #403 opened on Oct 21, 2016 by georgid Makam	idio alignment interfa	nce makam				Ļ	] <mark>6</mark>

# Keep track of versions

https://github.com/sertansenturk/ahenkidentifier/tags									
This repository Sea	Pull requests Issues Gist	🦨 ++ 💽 •							
📮 sertansenturk / ahe	nkidentifier	O Watch →1★ Star1% Fork0							
<> Code (!) Issues (3)	ן Pull requests וווי Projects וווי שואני אי Pulse וויא אין אין אין אין אין אין אין אין אין א	III Graphs							
Releases Tags									
on Jan 25	v1.5.2 -O- f174971 ≧ zip ≧ tar.gz ≧ Notes	Read release notes							
on Jun 1 2016	v1.5.0 … -O-fc14dfc 〗zip 〗tar.gz 〗Notes	Read release notes							
on Apr 19 2016	v1.4.0 … -0- 1d1b613 ≧ zip ≧ tar.gz ≧ Notes	Read release notes							
on Mar 16 2016	v1.3.0 … -O- c2a956f ≧ zip ≧ tar.gz ≧ Notes	Read release notes							
on Mar 16 2016	v1.2.0 -O- b6b9495 ≧ zip ≧ tar.gz ≧ Notes	Read release notes							
on Mar 4 2016	v1.1.1 -O- c782c5e ≧ zip ≧ tar.gz ≧ Notes	Read release notes							
on Mar 4 2016	v1.1	Read release notes							

### Keep track of versions

 Use git tags to mark the version that you used in a particular paper

 Do you want to try something new? make a branch and experiment without affecting your existing code

### Data

 Your code probably processes data. If you're distributing code, can you also distribute data?

- This can be more difficult:
  - Is the data yours to distribute? copyright/personally identifying?

Where do you host it?



- \* <u>zenodo.org</u>
- Designed for hosting public datasets
- Developed and maintained at CERN
- Hosted on the same storage infrastructure that stores data from the Large Hadron Collider
- Very generous storage limits (up to 50GB per dataset, email them if you need more!)

### zenodo

Q Upload

November 27, 2016

Available in

Publication date:

November 27, 2016

DOI 10.5281/zenodo.174616

### MTG/SymbTr-pdf: SymbTr-pdf v2.4.3

Sertan Senturk

This release contains the music scores in the pdf format of the SymbTr collection v2.4.3: https://github.com/MTG/SymbTr/tree/v2.4.3

#### Preview

SymbTr-pdf-v2.4.3.zip

#### ! The previewer is not showing all the files

#### MTG-SymbTr-pdf-1e8b36d

0	🗋 .gitignore	10 Bytes
0	AUTHORS	122 Bytes
0		201 Bytes
0	BEADME.md	787 Bytes
0	🗅 acemilahiduyekaldanma_dunyazekai_dede.pdf	163.4 kB
0	🗅 acemilahinimevsatcalabim_birhaci_bayram_veli.pdf	136.1 kB
0	🗅 acemkupeduyekzulfunuahmet_avni_konuk.pdf	35.9 kB
0	🗅 acemselamdevrikebirasik-i_gerhuseyin_fahreddin_dede.pdf	193.1 kB
0	🗅 acemseyirsofyan1erol_bingol.pdf	29.0 kB
0	🗅 acemseyirsofyan1sefik_gurmeric.pdf	65.4 kB
0	🗅 acemturkuduyekordunun_dereleriordu.pdf	139.2 kB
0	🗅 acemasiranaranagmeagiraksak1pdf	103.2 kB
0	🗅 acemasiranaranagmeaksak1pdf	99.1 kB
0	🗋 acemasiranaranagmecurcuna1pdf	102.3 kB
0	🗅 acemasiranaranagmesemai1pdf	100.7 kB
0	🗅 acemasiranaranagmesenginsemai1pdf	106.3 kB

Files (231.2 MB)		~
Name	Size	
MTG/SymbTr-pdf-v2.4.3.zip	231.2 MB	Preview & Download
md5:5bb477412b7d6ed816e527c17f9b067b 📀		

#### x

Š

Software Open Access

### DOI: Related identifiers: Supplement to: https://github.com/MTG/SymbTr-pdf/tree/v2.4.3 License (for files): C Other (Open) Share Cite as

Sertan Senturk. (2016, November 27). MTG/SymbTrpdf: SymbTr-pdf v2.4.3 (Version v2.4.3). Zenodo. http://doi.org/10.5281/zenodo.174616

GitHub

Start typing a citation style ...

#### Export

BibTeX CSL DataCite Dublin Core JSON JSON-LD MARCXML C Mendeley

### Zenodo + GitHib

 You can link a repository on GitHub to Zenodo so that every time you make a release you get a DOI to reference in your paper

## Data Backups

- Is your data backed up?
- \* Thirty-five minutes spent in Langley's Willowbrook Shopping Centre cost a Surrey woman much more than she had anticipated. Langley RCMP say that while she was shopping from I-1:35 p.m. last Monday, someone broke into her vehicle and stole a number of items, including a Mac iBook laptop containing the research she had compiled as she worked towards her PhD. "All that information was on that computer and she has no back-up file," said Langley RCMP spokesman Cpl. Brenda

Plumbley et al. 2012. Research reproducibility workshop

## Data Backups

- \* Dropbox
- \* External hard drive
- Department servers/backup infrastructure

## Licensing

I found this code online which does something I need to do. I can download the code, so I can use it, right?

gregversteeg / gaussianize					Watch      ▼	2 ★ Star	16	% Fork	2
<> Code	() Issues ()	្រា Pull requests 0	Projects 0 🗉 Wiki	-/- Pulse	III Graphs				
Transform data into normally distributed data.									
T commits I branch O r			♡ 0 releases	11	1 contributor		a <u>t</u> a N	ЛІТ	
		Tree: 7eb6aa9a70 - New pull request							
Tree: <b>7eb6aa</b>	a9a70 <del>-</del> New pu	Il request		Create new fi	le Upload files	Find file	Clone	or download	1-
Tree: 7eb6aa	a9a70 → New pu steeg Proofreading	ull request		Create new fi	le Upload files La	Find file	Clone 7eb6aa9	or download on Oct 5 201	1 -
Tree: 7eb6aa	a9a70 - New pu steeg Proofreading	ull request g. Added boxcox and brut	e force methods to gaussianiz	Create new fi	le Upload files La	Find file	Clone 7eb6aa9	or download on Oct 5 201 2 years ag	15 go
Tree: 7eb6aa	a9a70 - New pu steeg Proofreading	g. Added boxcox and brut Added boxcox and brut	e force methods to gaussianiz	Create new fi ze data. ze data.	le Upload files La	Find file	Clone 7eb6aa9	or download on Oct 5 201 2 years ag 2 years ag	15 go

# Licensing

📮 gregve	⊙ Watch ▼	5									
<> Code	() Issues 1	ື່ Pull requests 0	Projects 0	💷 Wiki	Insights						
License #1 (Closed alastair opened this issue on Feb 24, 2017 · 1 comment											
	alastair comment	ed on Feb 24, 2017				+ 💼 🥖	•				
	Hi, What is the license	e of this code?									
	gregversteeg cor	mmented on Feb 24, 201	17			Owner +	)				
	Thanks, I forgot. Added MIT license.										
	👍 1										
	🖉 🧝 gregverst	eeg closed this on Feb 2	24, 2017								

### Licensing

 No. Copyright law says that when you create something that no one can use or change it (without your permission)

gregversteeg / gaussianize					⊙ Watch -	2	🛨 Star	16	Υ <mark>γ</mark> Fork	2	
<> Code	! Issues 0	1 Pull requests 0	Projects 0	🗉 Wiki	Pulse	III Graphs					
Transform data into normally distributed data.											
T 8 T	To 8 commits			eleases	11	1 contributor		MIT د <u>أ</u> ه			
Branch: maste	er 🔻 New pull re	equest			Create new fi	le Upload file	s Fi	nd file	Clone	or download	1-
R gregvers	steeg committed o	n GitHub add license					Lates	st commi	23d0e4	6 19 days a	go
tests		Added boxcox and brut	te force methods to	gaussianiz	e data.					2 years a	go
LICENSE.	.txt	add license 19 days						19 days a	go		
🖹 gaussiani	ize.py	Added boxcox and brute force methods to gaussianize data. 2 years						2 years a	go		
🖹 readme.n	readme.md Proofreading. 2 years							2 years a	go		

# Open licenses

- Two main styles of open license:
  - GPL: "Use this code, but if you change it or integrate it into another system, your changes and that system also need to be available under this same open license"
  - BSD/MIT/Apache: "Do whatever you want. Really. Just say that I made it"
- Check with the copyright holder of your code (It might be the university)

### How do you add a licence?

- Make a file called LICENCE or COPYING containing the contents of the licence
- GitHub will do this for you when you make a repository. It will ask you what licence you want and will automatically add the file

### Reproducible research

- Sertan Şentürk (MTG PhD student) gave an example in a talk —3 weeks to reproduce an experiment with an external researcher (<u>https://zenodo.org/record/255537</u>)
- \* Do you know what the dependencies are? How to format the input data? How to execute the software? Does your code still work?
- \* This is not just for others. For you too
- "The first person to benefit by making your work reproducible is you!" - Sertan Şentürk

### Reproducible research

 It's my belief that following these recommendations will make it easier to have research that can be verified and extended

 An entirely reproducible paper. Can someone else (you) run all of your experiments, create all tables, graphs, and generate your paper/thesis by running just I command? (Victoria Stodden, University of Illinois at Urbana-Champaign)

# Wrapping up

\* This is a lot of work. Do I have to do it?

- You're saving time
- \* You'll become a better programmer (easier to get a job)
- \* You'll actually finish your PhD on time!

### Thanks!







This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowsa-Curie grant agreement No. 765068.