

Reproducible research and reproducible analyses: literate programming for everyone

Adam J. Richards

Centre national de la recherche scientifique (CNRS)
(French National Center for Scientific Research)
Station d'Ecologie Expérimentale du CNRS à Moulis

Last updated: November 23, 2012

Why are scientific studies so difficult to reproduce?

- Publication/Experimental bias
- Rewards for 'positive results'
- Programming errors or data manipulation mistakes
- Poorly selected statistical tests
- Multiple testing, multiple looks at the data, multiple statistical analyses
- Lack of easy-to-use tools



Ideas adopted from a presentation by Jim Berger
Image taken from <http://wanderingfrance.com>



- Bayer Healthcare reviewed 67 in-house attempts at replicating the findings in published research.
 - Less than 1/4 were viewed as essentially replicated.
 - Over 2/3 had major inconsistencies → project termination.
- *Why Most Published Research Findings Are False*
- A presentation by Jim Berger

J.P. Ioannidis. Why Most Published Research Findings Are False **PLoS Med.** 2005 August; 2(8): e124. [URL]

Why bother making our work reproducible?

Nature Medicine - 12, 1294 - 1300 (2006)

Genomic signatures to guide the use of chemotherapeutics

Anil Potti^{1,2}, Holly K Dressman^{1,3}, Andrea Bild^{1,3}, Richard F Riedel^{1,2}, Gina Chan⁴, Robyn Sayer⁴,
 Janiel Cragun⁴, Hope Cottrill⁴, Michael J Kelley², Rebecca Petersen⁵, David Harpole⁵, Jeffrey Marks⁵,
 Andrew Berchuck^{1,6}, Geoffrey S Ginsburg^{1,2}, Phillip Febbo^{1,2,3}, Johnathan Lancaster⁴ & Joseph R Nevins^{1,2,3}

Dr. Baggerly and Dr. Coombes found errors almost immediately. Some seemed careless – moving a row or a column over by one in a giant spreadsheet – while others seemed inexplicable. The Duke team shrugged them off as “clerical errors.”

And the Duke researchers continued to publish papers on their genomic signatures in prestigious journals. Meanwhile, they started three trials using the work to decide which drugs to give patients.

source: <http://www.nytimes.com/2011/07/08/health/research/08genes.html>

Again why bother?

- Retracted: 07 January 2011
- Streamlines the review of manuscripts and grant proposals
- Like writing good code—it ultimately saves time
- Promotes sharing
- Promotes **good practices**

Multiple testing

The tradition in epidemiology is to ignore multiple testing

Over-reliance on the use of p -values

The tradition in psychology is to ignore optional stopping; if the p -value is close to significant then get one more data point (with no adjustment).

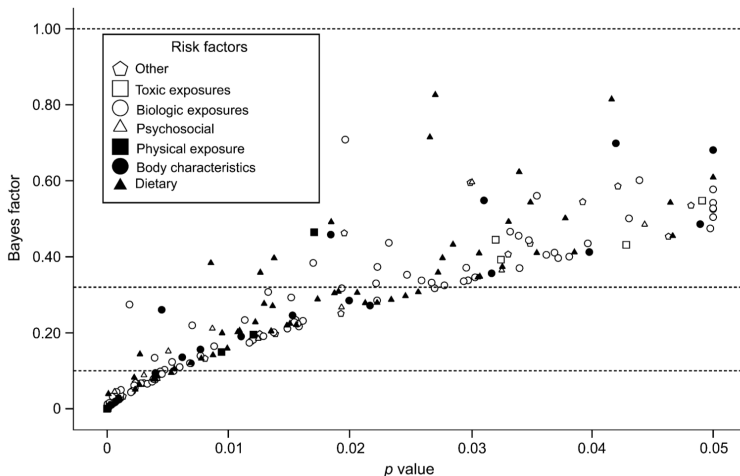


FIGURE 1. Estimated Bayes factors for 272 epidemiologic studies with formally statistically significant results. The Bayes factor is plotted against the observed p value in each study. Shown are calculations assuming θ_A of 0.50 (relative risk = 1.65). The dashed lines correspond to threshold values (1.00, 0.32, 0.10) separating different Bayes factor categories.

J P A Ioannidis *Am. J. Epidemiol.* (2008) 168 (4): 374-383. [URL]

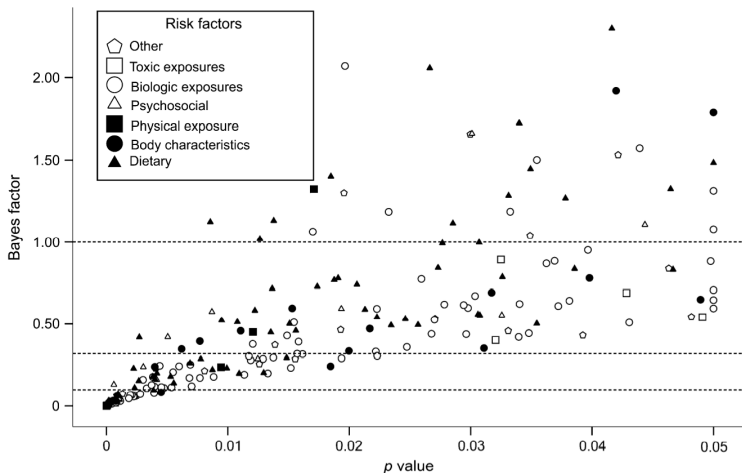


FIGURE 2. Estimated Bayes factors for 272 epidemiologic studies with formally statistically significant results. The Bayes factor is plotted against the observed p value in each study. Shown are calculations assuming θ_A of 1.50 (relative risk = 4.48). The dashed lines correspond to threshold values (1.00, 0.32, 0.10) separating different Bayes factor categories.

J P A Ioannidis *Am. J. Epidemiol.* (2008) 168 (4): 374-383. [[URL](#)]

Reproducible Research

Reproducible research consists two main components:

- ① Reproducible experiment
- ② Reproducible analysis

Lab/Experiment

Data Sharing	Instruments, samples, buffers	Supplemental Methods
Audit Trail	Version control and backup	Git, Mercurial and SVN
Documentation	All details required to reproduce	Electronic lab notebook

Data Analysis

Data Sharing	raw, standards	public repositories, URL
Audit Trail	Version control and backup	Git, Mercurial and SVN
Documentation	All details required to reproduce	Literate Programming

Data Sharing

- 1 Must include the raw data and the appropriate meta data
- 2 Must be free and publicly available
- 3 Use data standards when available (eg. MIAME or MINSEQE)

Repositories

- Microarray - [Gene Expression Omnibus](#) - NCBI
- Transcriptomics data - [ArrayExpress](#) - EBI
- Sequencing data - [Read Archive](#) - NCBI
- Sequencing data - [Trace Archive](#) - NCBI
- Metagenomics - [Metagenomics](#) - EBI
- Genome-phenome - [European Genome-phenome Archive](#) - EBI

Audit Trail

What kinds of systems are available?

- **Good** - The cloud (Dropbox, Google Drive)
- **Better** - Version control systems (SVN, Git and Mercurial*)
- **Best** - Version control systems on the cloud ([GitHub](#), [Bitbucket](#))

Good practices

- 1 Use a system that documents both data and process
- 2 Use the machine readable CSV format
- 3 Never embed data manipulation and statistical tests
- 4 Use R, Python or another freely available software to read and process raw data—ideally to produce reports complete with code, results and prose.

* For Windows users look at [TortoiseHg](#).

Create a wiki and track changes to it through a repository

- 1 Create a repository called 'labwiki' on bitbucket or on another hosted site
- 2 Download and setup wiki

```
$ hg clone https://<USERNAME>/bitbucket.org/labwiki
$ cd labwiki
$ wget pmwiki.org/pub/pmwiki/pmwiki-latest.tgz
$ tar zxvf pmwiki-latest.tgz
$ cd pmwiki-*
```

- 3 Create a config.php file (use the one from the site)
- 4 add it to your repository

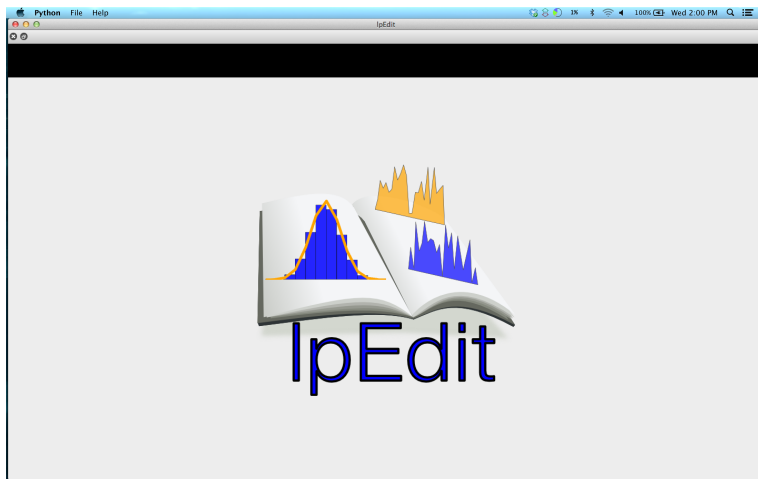
```
$ hg add *
$ hg commit -m "my first commit" --user <ajrichards>
$ hg push
```

- 5 copy it to a public_html directory (FileZilla if remote)

```
$ cp -r blah /public_html
```

IpEdit

:



About IpEdit

IpEdit has two parts:

- ① An editor built specifically for literate programming.
- ② A webpage that contains the editor documentation, examples and acts a general resource for reproducible research

IpEdit is

- ① cross-platform
- ② open source and freely available
- ③ Easy to use
- ④ based on QScintilla and written using PyQt4
- ⑤ working with both R and Python languages

Literate Programming

The concept was introduced by Donald Knuth in the 1970's. The idea is that we should be able to read and write code as if it were an expression of logic. Thus, the data, code and text used to tell the story must be presented as one single report.

A good literate programming tool should be...

- 1 Easy to use
- 2 Produce attractive output formats
- 3 Be capable of producing PDF, HTML and presentations
- 4 Should not be editor specific
- 5 Free and available to everyone
- 6 Well documented

syntax based on noweb

This is text that goes around the code

```
<<label=chunk1>>=  
your code  
@
```

This is more text that goes around the code

Basic Sweave program

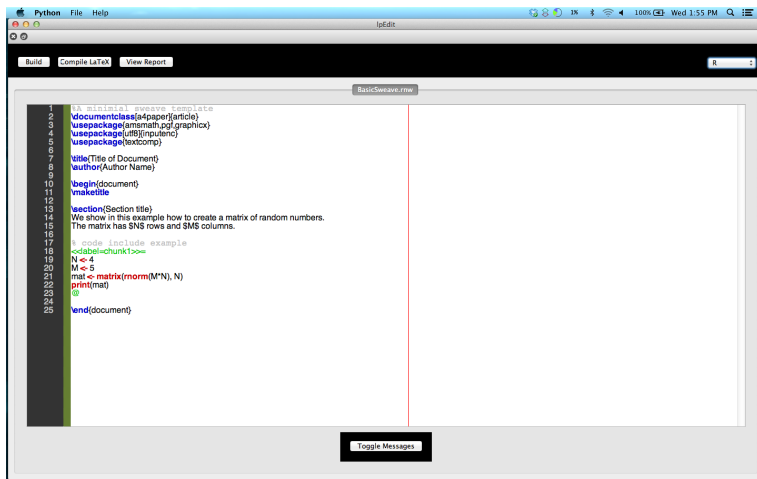
```
\documentclass[a4paper]{article}
\usepackage{amsmath,pgf,graphicx,textcomp}
\usepackage[utf8]{inputenc}
\title{Title of Document}
\author{Author Name}

\begin{document}
\maketitle

\section{Section title}
This matrix has  $N$  rows and  $M$  columns.

<<label=chunk1>>=
N <- 4
M <- 5
mat <- matrix(rnorm(M*N), N)
print(mat)
@
\end{document}
```

lpEdit - DEME



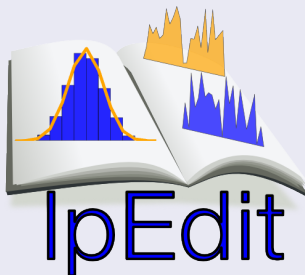
The screenshot shows the lpEdit application window. The title bar indicates it's a Python application. The menu bar includes Python, File, and Help. The status bar at the top shows various system icons, including network, battery, and time (Wed 1:55 PM). The main window has a toolbar with buttons for Build, Compile LaTeX, and View Report. A tab labeled 'Basic5weave.rnw' is active. The editor displays a LaTeX document with the following content:

```
1 \documentclass[12pt]{article}
2 \documentclass[a4paper]{article}
3 \usepackage{amsmath,pdf,graphicx}
4 \usepackage{utf8}(inputenc)
5 \usepackage{textcomp}
6
7 \title{Title of Document}
8 \author{Author Name}
9
10 \begin{document}
11 \maketitle
12
13 \section{Section title}
14 We show in this example how to create a matrix of random numbers.
15 The matrix has  $N$  rows and  $M$  columns.
16
17 % code include example
18 <label=chunk1>
19 N <- 4
20 M <- 5
21 mat <- matrix(rnorm(M*N), N)
22 print(mat)
23 @
24
25 \end{document}
```

At the bottom of the window, there is a 'Toggle Messages' button.

In the future...

- we may include additional languages
- we will have many more examples
- will produce HTML output at the click of a button
- will be working with restructured text markup (reST)
- will have more customization options
- code folding, inline spellcheck



In an ideal world...

Acknowledgments

⋮

Reproducible Analysis

All raw data and associated metadata described in this manuscript were deposited in the public repository *<http://myrepository.org>* and are available without restriction (for non-commercial purposes) to the general public.

All data manipulations and statistical analyses presented here are documented using literate programming and the complete bundle containing example data, code and the final report is available as supplemental methods.

Importantly, these groups have generously provided funding.

- Duke University - Duke Cancer Institute (DCI)
- Centre national de la recherche scientifique (CNRS)

These people have contributed to the project.

- Kouros Owzar (Duke University)
- Andrzej Kosinski (Duke University)