# Reproducible research in data science - a bioinformatics primer

Mikhail Dozmorov, Ph.D.
Department of Biostatistics, VCU
mikhail.dozmorov@vcuhealth.org
01/29/2016

---

# Reproducible research in science

- Science is the systematic enterprise of gathering knowledge about the universe and organizing and condensing that knowledge into testable laws and theories.

- The success and credibility of science are anchored in the willingness of scientists to expose their ideas and results to independent testing and replication by other scientists.

http://www.aps.org/policy/statements/99_6.cfm

---

# What is reproducible research?

- Reproducibility
- Replicability
- Repeatability
- Reliability
- Robustness
- Generalizability

TRUTH  Open  TRANSPARENCY

Steve Goodman, Stanford, March 18, 2015

---

# What is reproducible research?

Reproducible research is the ultimate standard for strengthening scientific evidence by independent:

- Investigators
- Data
- Analytical methods
- Laboratories
- Instruments

Replication is particularly important in studies that can impact broad policy or regulatory decisions

http://www.nature.com/news/robust-research-institutions-must-do-their-part-for-reproducibility-1.18259

---

# Why do we care?

- High-throughput data-generating technologies are increasingly used to make clinical recommendations and treatment decisions
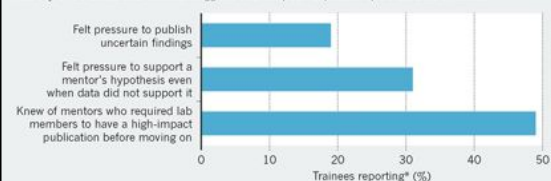- A problem may be overlooked .. Published .. Get in clinical trials



http://retractionwatch.com/2011/05/04/the-importance-of-being-reproducible-keith-baggerly-tells-the-anil-potti-story/

---

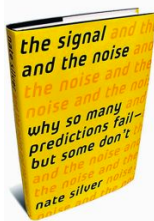# Why reproducible research is questioned now?

- Publish or perish



http://www.nature.com/news/robust-research-institutions-must-do-their-part-for-reproducibility-1.18259

## Why irreproducibility problems arise?

- Humans are good at recognizing patterns

Human beings do not have very many natural defenses. We are not all that fast, and we are not all that strong. We do not have claws or fangs or body armor. We cannot spit venom. We cannot camouflage ourselves. And we cannot fly. Instead, we survive by means of our wits. Our minds are quick. **We are wired to detect patterns** and respond to opportunities and threats without much hesitation.
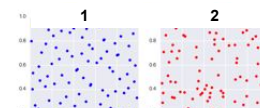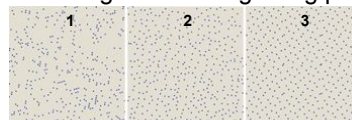
- Nate Silver

*the signal and the noise — why so many predictions fail- but some don't — nate silver*

http://www.amazon.com/gp/product/0143125087?redirect=true&ref_=s9_simh_gw_g14_i1_r
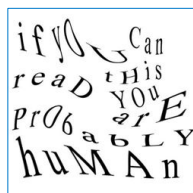
---

## Why irreproducibility problems arise?

- Humans are good at recognizing patterns



http://bit-player.org/2011/a-slight-discrepancy    https://sobol.readthedocs.org/en/latest/

---

## Why irreproducibility problems arise?

- Humans are good at recognizing patterns



overlooks    inquiry

http://neuralnetworksanddeeplearning.com/chap1.html
http://www.npr.org/templates/story/story.php?storyId=130594039

---

## Why irreproducibility problems arise?

- Humans are good at recognizing patterns

... improved technology did not cover for the lack of theoretical understanding about the economy, it only gave economists faster and more elaborate ways to mistake noise for a signal.
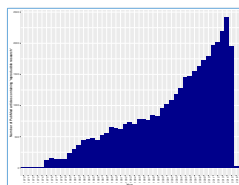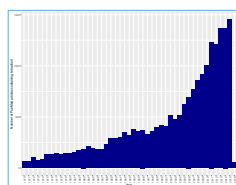
- Nate Silver

*the signal and the noise — why so many predictions fail- but some don't — nate silver*

**Essay**

*Open access, freely available online*

### Why Most Published Research Findings Are False

John P.A. Ioannidis

http://www.amazon.com/gp/product/0143125087?redirect=true&ref_=s9_simh_gw_g14_i1_r
http://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.0020124

---

## PubMed on "Reproducible research" vs. "Retraction"

"Reproducible research"      "Retraction"



Number of publications per year, from 1970 to January 2016

Retraction Watch

Tracking retractions as a window into the scientific process

http://retractionwatch.com/

---

## Reproducibility initiatives

**RESEARCH ARTICLE**

PSYCHOLOGY

### Estimating the reproducibility of psychological science

Open Science Collaboration*†

Reproducibility is a defining feature of science, but the extent to which it characterizes current research is unknown. We conducted replications of 100 experimental and correlational studies published in three psychology journals using high-powered designs and original

**Reproducing 100 psychological studies**

significant results. Thirty-six percent of replications had statistically significant results; 47% of original effect sizes were in the 95% confidence interval of the replication effect size; 39% of effects were subjectively rated to have replicated the original result; and if no bias in original results is assumed, combining original and replication results left 68% with statistically significant effects. Correlational tests suggest that replication success was better predicted by the strength of original evidence than by characteristics of the original and replication teams.

http://science.sciencemag.org/content/sci/349/6251/aac4716.full.pdf

Validation.

Reproducibility Project:
Cancer Biology

About

**Reproducing 50 (now 37) cancer studies**

http://validation.scienceexchange.com/#/cancer-biology

## Reproducibility initiatives

nature.com

**Enhancing reproducibility**
New reporting standards for Nature journal authors are intended to improve transparency and reproducibility.

https://www.ncbi.nlm.nih.gov/pubmed/23762900

NEWS

**REPRODUCIBLE RESEARCH**
ADDRESSING THE NEED FOR DATA AND CODE SHARING IN COMPUTATIONAL SCIENCE
*By the Yale Law School Roundtable on Data and Code Sharing*

https://web.stanford.edu/~vcs/papers/RoundtableDeclaration2010.pdf

NIH plans to enhance reproducibility

**Francis S. Collins** and **Lawrence A. Tabak** discuss initiatives that the US National Institutes of Health is exploring to restore the self-correcting nature of preclinical research.

https://www.ncbi.nlm.nih.gov/pubmed/24482835

## Reproducibility initiatives

C✲S    CENTER FOR OPEN SCIENCE

https://cos.io/

**Open Science Framework**
A scholarly commons to connect the entire research cycle

https://osf.io/

Order experiments from the world's best labs

https://www.scienceexchange.com/

## WHAT CAN WE DO TO ENHANCE REPRODUCIBILITY?

## Steps in reproducible research

The most important tool is the mindset, when starting, that the end product will be reproducible.
– Keith Baggerly

- Can my conclusions be reproduced by others?
- Data availability
- Detailed Methods description
- Software availability

- Empirical reproducibility
- Computational reproducibility
- Statistical reproducibility

http://kbroman.org/Tools4RR/

## Common approach: write report around results

**Point and click approach**
- Use MS Excel for data entry/cleaning/preparation, and possibly statistical analysis;
- Copy/paste to/from other programs.

Zeeberg BR et al. **Mistaken identifiers: gene name errors can be introduced inadvertently when using Excel in bioinformatics** *BMC Bioinformatics*. 2004
https://www.ncbi.nlm.nih.gov/pubmed/15214961

**Problems**
- With point-and-click, there's no way to record/save the steps that generated the results;
- Data files are kept separately from the analysis code, and from reports;
- After modifications of one of the files involved, it becomes unclear which version corresponds exactly to the reported results;
- Every time something changes, you have to regenerate the figures/results/reports by hand – very time consuming.

## Better approach: write report that generates results

- Everything automated via code;
- Most raw data is attached to the code;
- Any changes in code should be version controlled;
- The full report should be self-sufficient and reproducible with a single command.

## EVERYTHING AUTOMATED

---

## Automating everything, or, why we love R

- R is an open source programming language – removes error-prone point-and-click;
- R is free to run, study, change and improve;
- R runs on Windows, Mac, and Linux;
- R has publication quality graphing capabilities;
- Extensible with a very large collection of actively developing packages;
- Excellent report-creating capabilities.

http://rpubs.com/ideltoro/147317

**The New York Times**

"**R is really important to the point that it's hard to overvalue it,**" said Daryl Pregibon, a research scientist at Google, which uses the software widely. "**It allows statisticians to do very intricate and complicated analyses without knowing the blood and guts of computing systems.**"
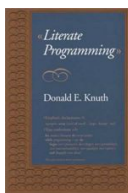
http://www.nytimes.com/2009/01/07/technology/business-computing/07program.html

---

## R is reimagined with RStudio

Code editor, multiple files

Console to execute the code

Variables, data, history of commands

Plots, help, files, packages

---

## CODE ATTACHED TO THE REPORT

---

## Literate programming, or self-documenting code

- A report containing a stream of text and code chunks;
- Each code chunk loads data, computes results, shows figures;
- Each text chunk explains how the code chunks work;
- The resulting report is human- and machine readable.

*Literate Programming*
Donald E. Knuth

https://en.wikipedia.org/wiki/Literate_programming

---

## Evolution of literate programming

- **LaTeX** – document markup language;
- **Sweave** – integrates LaTeX markup and R code formatting;
- **R markdown (knitr)** – most natural way of combining text and code.

| LaTeX | R Markdown |
|---|---|
| \textbf{word} | **word** |

http://arxiv.org/abs/1402.1894

## Markdown basics

Markdown Quick Reference ▾

**Emphasis**

*italic*  **bold**

_italic_  __bold__

**Headers**

# Header 1

## Header 2

### Header 3

**Lists**

Unordered List

* Item 1
* Item 2
    + Item 2a
    + Item 2b

Ordered List

1. Item 1
2. Item 2
3. Item 3
    + Item 3a
    + Item 3b

```
## Getting started
To work with R Markdown, if necessary:

* Install [R](http://www.r-project.org/)
* Install the lastest version of [RStudio](http://rstudio.org/
download/) (at time of posting, this is 0.96)
* Install the latest version of the `knitr` package:
`install.packages("knitr")`

To run the basic working example that produced this blog post:

* Open R Studio, and go to File - New - R Markdown
* If necessary install `ggplot2` and `lattice` packages:
`install.packages("ggplot2"); install.packages("lattice")`
* Paste in the contents of this gist (which contains the R
Markdown file used to produce this post) and save the file with
an `.rmd` extension
* Click Knit HTML
```

**Getting started**

To work with R Markdown, if necessary:

- Install R
- Install the lastest version of RStudio (at time of posting, this is 0.96)
- Install the latest version of the `knitr` package: `install.packages("knitr")`

To run the basic working example that produced this blog post:

- Open R Studio, and go to File - New - R Markdown
- If necessary install `ggplot2` and `lattice` packages: `install.packages("ggplot2"); install.packages("lattice")`
- Paste in the contents of this gist (which contains the R Markdown file used to produce this post) and save the file with an `.rmd` extension
- Click Knit HTML

---

## Literate programming with knitr

- Mix markdown with code

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```{r}
summary(cars)
```

You can also embed plots, for example:

```{r, echo=FALSE}
plot(cars)
```

- And "knit" a report with one button

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(cars)

##      speed           dist
## Min.   : 4.0    Min.   :  2.00
## 1st Qu.:12.0    1st Qu.: 26.00
## Median :15.0    Median : 36.00
## Mean   :15.4    Mean   : 42.98
## 3rd Qu.:19.0    3rd Qu.: 56.00
## Max.   :25.0    Max.   :120.00
```

---

# VERSION CONTROL OF CHANGES

---

## Why we need version control?



http://www.phdcomics.com/comics/archive.php?comicid=1531

---

## Version control – what and when did you do

- **Git** and GitHub – version control system;
- Each project stored in its own repository;
- History of changes – track what you did;
- Ability to go back if something breaks;
- Branch out, go creative, then merge or revert the changes;
- Collaborate through merging changes from multiple people.
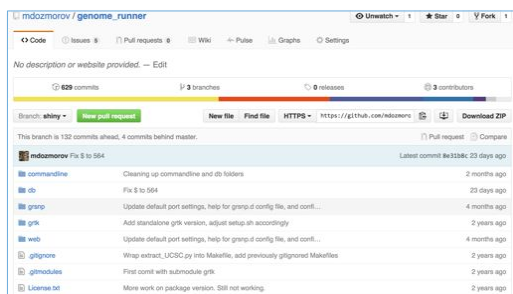
https://github.com/

---

## Git basics

- Git is a command line tool;
- GitHub is a web-based storage for your project repositories;
- **Git add** – add a file to version control system;
- **Git commit** – make a snapshot of current changes;
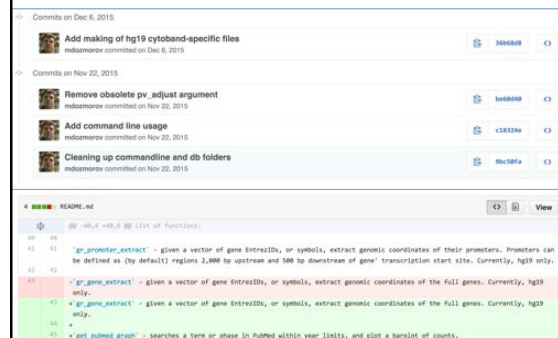- **Git push/pull** – send/get changes to/from GitHub.

https://github.com/

# Repository on GitHub



# History of commits



# Publishing with Git

**B-Cell and Monocyte Contribution to Systemic Lupus Erythematosus Identified by Cell-Type-Specific Differential Expression Analysis in RNA-Seq Data**

Libertas Academica

Mikhail G. Dozmorov[1], Nicolas Dominguez[2], Krista Bean[2], Susan R. Macwana[2], Virginia Roberts[2], Edmund Glass[1], Judith A. James[2] and Joel M. Guthridge[2]

**Implementation and availability.** All RNA-seq data processing steps were performed in CentOS 6.6 high-performance cluster computing environment. All analyses were conducted in R/Bioconductor environment v 3.2.0.[36,37] All analytical scripts are available at https://github.com/mdozmorov/deconvolution.

https://www.ncbi.nlm.nih.gov/pubmed/26512198

# LEARN MORE

# Reproducible research made simple

OPEN ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

**Editorial**

## Ten Simple Rules for Reproducible Computational Research

Geir Kjetil Sandve[1,2*], Anton Nekrutenko[3], James Taylor[4], Eivind Hovig[1,5,6]

https://www.ncbi.nlm.nih.gov/pubmed/24204232

## Best Practices for Scientific Computing

Greg Wilson [*], D.A. Aruliah [†], C. Titus Brown [‡], Neil P. Chue Hong [§], Matt Davis [¶], Richard T. Guy [∥], Steven H.D. Haddock [**], Katy Huff [††], Ian M. Mitchell [‡‡], Mark D. Plumbley [§§], Ben Waugh [¶¶], Ethan P. White [***], Paul Wilson [†††]

https://www.ncbi.nlm.nih.gov/pubmed/24415924

# Reproducible research in data science - a bioinformatics primer

## Thank you

https://github.com/mdozmorov/presentations

Mikhail Dozmorov, Ph.D.
Department of Biostatistics, VCU
mikhail.dozmorov@vcuhealth.org
01/29/2016