



Connecting R/R Markdown and Microsoft Word using StatTag for Collaborative Reproducibility

R-Ladies

September 26, 2019

Evanston, Illinois

Leah J. Welty, PhD

Disclosures

- Joint work with:
 - Luke V. Rasmussen, Lead Software Developer
 - Abigail S. Baldridge, Senior Statistical Analyst
 - Eric W. Whitley, Software Developer
- This project was supported, in part, by the National Institutes of Health's [National Center for Advancing Translational Sciences](#), Grant Number **UL1TR001422**. *The content is solely the responsibility of the developers and does not necessarily represent the official views of the National Institutes of Health.*





Reproducible Research: Background

Reproducible Research: Background

Origins lie in the inconvenience of *irreproducible* research

MAKING SCIENTIFIC COMPUTATIONS REPRODUCIBLE

To verify a research paper's computational results, readers typically have to recreate them from scratch. ReDoc is a simple software filing system for authors that lets readers easily reproduce computational results using standardized rules and commands.

“In the mid-1980s, we realized that our laboratory’s researchers often had difficulty reproducing their own computations without considerable agony.”

In the mid-1980s, we realized that our laboratory’s researchers often had difficulty reproducing their own computations without considerable agony. We also noticed that junior students, who typically build on the work of more advanced students, frequently spent a great deal of time and effort just to reproduce their colleagues’ computational results.

Reproducing computational research poses challenges in many environments. Indeed, the problem occurs wherever people use the traditional methods of scientific publication to describe computational research. For example, in a traditional article, the author simply outlines the relevant computations—the limitations of a paper medium prohibit complete documentation, which would ideally include experimental data, parameter values, and the author’s programs. Readers who wish to use and verify the work must reimplement it, which is often a painful process. Even if readers have access to the author’s source files (a feasible assumption given re-

cent progress in electronic publishing) they can only recompute the results by invoking the various programs exactly as the author invoked them; such information is something that is usually undocumented and difficult to reconstruct.

To address these problems, we developed ReDoc, a system for reproducing scientific computations in electronic documents. Since implementing it in the early 1990s, ReDoc has become our principal means for organizing and transferring our laboratory’s scientific computational research.

ReDocs are best defined operationally: After an author completes a ReDoc, readers can destroy all existing results—principally the illustrations—and rebuild them using the author’s underlying programs and raw data. Using ReDoc, authors describe their computations and preserve their details in fully functional examples. Authors can also test their archived research software by occasionally removing and regenerating the document’s results using ReDoc’s standardized interface commands. ReDoc also lets authors develop automatic scripts to verify any document’s completeness and reproducibility before its publication. Scientific journal publishers might also use ReDoc in the referee process to test the reproducibility of illustrations.

ReDoc benefits readers in several ways. Just as a driver wants to find the brake pedal in the

1521-9615/00/\$10.00 © 2000 IEEE

MATTHIAS SCHWAB, MARTIN KARRENBACH,
AND JON CLAERBOUT
Stanford University

Reproducible Research: Dynamic Documents

Synonymous with Dynamic Documents

Sweave: Dynamic Generation of Statistical Reports Using Literate Data Analysis

Friedrich Leisch

Institut für Statistik und Wahrscheinlichkeitstheorie, Technische Universität Wien, Wiedner Hauptstraße 8-10/1071, A-1040 Wien, Austria

Abstract. Sweave combines typesetting with L^AT_EX and data analysis with S into integrated statistical documents. When run through R or Splus, all data analysis output (tables, graphs, ...) is created on the fly and inserted into a final L^AT_EX document. Options control which parts of the original S code are shown to or hidden from the reader, respectively. Many S users are also L^AT_EX users, hence no new software has to be learned. The report can be automatically updated if data or analysis change, which allows for truly reproducible research.

Keywords. R, S, literate statistical practice, integrated statistical documents, reproducible research

1 Introduction

The traditional way of writing a report as part of a statistical data analysis project uses two separate steps: First, the data are analyzed using one's favorite statistical software package, and afterwards the results of the analysis (numbers, graphs, ...) are used as the basis for a written report. In larger projects the two steps may be repeated alternately, but the basic procedure remains the same. Many statistical software packages try to support this process by generating pre-formatted tables and graphics that can easily be integrated into a final report using copy-and-paste from the data analysis system to the word processor. The basic paradigm is to write the report around the results of the analysis.

Another approach for integration of data analysis and document writing is to embed the analysis itself into the document, which reverses the traditional paradigm. Over the last decade a number of systems have been developed that integrate analysis and documentation and allow for *literate statistical practice*, see Rossini (2001) for a survey.

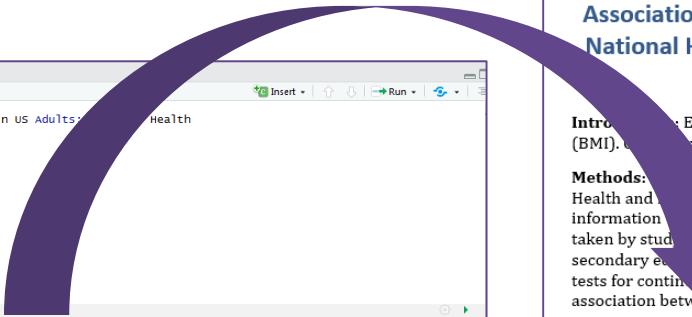
This new paradigm is probably most popular for creation of dynamic web pages and offers completely new possibilities for teaching statistics and delivering statistical methodology over the Internet. E.g., the ExploRe system (Härdle et al., 1999) provides means to embed statistical quantlets in web pages or electronic books to create interactive documents with direct access to a statistical data analysis package. Another example for a dynamic statistical analysis on a web page is given in Temple Lang (2001), by embedding R into netscape as a plugin. Report rendering is performed using XML and XSL.

We introduce a new system, called Sweave, which combines ideas from both worlds described above using literate programming tools. The purpose is to create dynamic reports, which can be updated automatically if data or analysis change, while using standard tools for both data analysis and word processing. Sweave is written in the S language, either the open source R (<http://www.R-project.org>) or the commercial Splus (<http://www.insightful.com>) can be used for statistical data analysis.

Sweave introduced in 2002.
Reproducibility becomes synonymous with using dynamic documents that combine manuscript/report with code and data (e.g. Sweave, knitR).

Dynamic Documents: A Cornerstone of Reproducibility

Why we love them



```
② NHANES Example.Rmd
1 --- 
2 title: 'Association of Education with Anthropometrics in US Adults: National Health and Nutrition Examination Study 2013-2014'
3 author: 'John Doe'
4 geometry: width=75in
5 output: html_document
6 sansfont: calibri Light
7 fontsize: 11pt
8 ---
9 
10 <style type="text/css">
11   h1.title {
12     font-size: 11pt;
13     font-weight: bold;
14   }
15 </style>
16 
17 ---
18 ```{r setup, include=FALSE}
19 knitr::opts_chunk$set(echo = FALSE)
20 setwd("R:/NUCATS/NUCATS_Shared/BERDShared/Analysis Manager/Data and Programs/Reproducible Research Class/R Markdown/windows")
21 analysis<-read.csv("Analysis.csv")
22 attach(analysis)
23 library(tableone)
24 library(knitr)
25 options(digits=2)
26 
27 \usepackage[utf8]
28 
29 **Introduction:** Education level has been shown to be associated with body mass index (BMI). Gender, race, marital status and metabolic characteristics may alter this association.
30 
31 **Methods:** This study included adult ( $\geq 30$  years) participants from the 2013-2014 National Health and Nutrition Examination Study (NHANES). Education and demographic information were assessed by questionnaire and anthropometric measurements were taken by study personnel. Education was dichotomized based on matriculation to post-secondary education. Associations were estimated using T-tests or Wilcoxon rank sum tests for continuous data and chi-squared tests for categorical data. We examined the association between BMI and education level using multivariate linear regression.
32 
33 ```{r, echo=FALSE}
34 ps<-table(analysePostSecondary)
35 percents<-100*prop.table(ps)
36 model <- lm(BMXBMI ~ Postsecondary, data = analysis)
37 BetaUnivariable <- summary(model)$coefficients[2, 1]
38 LBUnivariable <- (summary(model)$coefficients[2, 1]) - 1.96 * (summary(model)$coefficients[2, 2])
39 UBUnivariable <- (summary(model)$coefficients[2, 1]) + 1.96 * (summary(model)$coefficients[2, 2])
40 myVars <- c("Gender", "Race", "RIDAGEYR", "Married", "BMXBMI", "LBXTC", "LBXGLU")
41 Tableone <- Createableone(data = analysis, vars=myVars, strata= "Postsecondary")
```

Association of Education with Anthropometrics in US Adults: National Health and Nutrition Examination Study 2013-2014

Introduction: Education level has been shown to be associated with body mass index (BMI). Gender, race, marital status and metabolic characteristics may alter this association.

Methods: This study included adult (≥ 30 years) participants from the 2013-2014 National Health and Nutrition Examination Study (NHANES). Education and demographic information were assessed by questionnaire and anthropometric measurements were taken by study personnel. Education was dichotomized based on matriculation to post-secondary education. Associations were estimated using T-tests or Wilcoxon rank sum tests for continuous data and chi-squared tests for categorical data. We examined the association between BMI and education level using multivariate linear regression.

Results: Among 4808 participants, 2649 (55.1%) self-reported any post-secondary education. Post-secondary education was associated with lower BMI (Beta: -0.62, 95% CI: -1.03 to -0.21). After adjusting for gender, race, age, marital status, fasting glucose and total cholesterol, post-secondary education was no longer significantly associated with BMI (Beta: -0.19, 95% CI: -0.79 to 0.41).

Table 1. Association of Education with Participant Characteristics among 2013-2014 NHANES Participants

	No Postsecondary	Postsecondary	p
n	2159	2649	
Gender = Male (%)	1079 (50.0)	1208 (45.6)	0.003
Race (%)			<0.001
Mexican American	456 (21.1)	173 (6.5)	
Non-Hispanic Asia	161 (7.5)	411 (15.5)	

- Numbers in document are updated automatically when data or models change.
- **Eliminates copying and pasting output.**
- Provides a link between a number in a manuscript and its provenance.

Reproducible Research: Sweave to R Markdown

Tools evolve: Sweave, knitR, R Markdown

```
\documentclass[a4paper]{article}

\title{Sweave Example 1}
\author{Friedrich Leisch}

\begin{document}

\maketitle

In this example we embed parts
\texttt{\{kruskal.test\} help pag}

<<>>=
data(airquality, package="data"
library("stats")
kruskal.test(Ozone ~ Month, da
@ which shows that the location distribution varies significantly
include a boxplot of the data:

\begin{center}
<<fig=TRUE,echo=FALSE>>=
library("graphics")
boxplot(Ozone ~ Month, data =
@ \end{center}

\end{document}
```

The screenshot shows the RStudio interface with two panes. The left pane displays the R Markdown source code, and the right pane shows the rendered output.

Left Pane (Source Code):

```
1 ---  
2 title: "R Markdown Example"  
3 author: "Leah Welty"  
4 date: "April 6, 2016"  
5 output: word_document  
6 ---  
7  
8 ````{r setup, include=FALSE}  
9 knitr::opts_chunk$set(echo = TRUE)  
10 ````  
11 you can use R Markdown from within RStudio. You write in a simple text editor, using the (fairly simple) Markdown language to indicate *italics* or **bold**. You can embed 'chunks' of R code and output them in the document.  
12  
13 For example, if I want to see a summary of the *cars* dataset that comes standard with R, I can insert R code that produces this:  
14  
15 ````{r cars}  
16 summary(cars)  
17 ````  
18  
19 I can also embed results directly in the text. For example, the median speed is `r mean(cars$speed)`.  
20  
21 That's pretty nice, because if I change something about the data, then that number can be automatically updated. This is how I'm changing the data:  
22  
23 ````{r newmean}  
24 cars$speed[1]  
25 cars$speed[1] <- 10  
26 ````  
27 So now if I generate the mean speed, it is `r mean(cars$speed)`.  
28  
29 You can also include plots, and make tables using R Markdown.  
30  
31 R Markdown will take your plain text file and at the touch of a button, insert all the R output then turn it into HTML, PDF, or MS Word. Pretty cool ... except ...  
32  
33 What happens when you send the word document to a collaborator, and they mark it up in track changes? [Hint: You end up abandoning R Markdown, or some unlucky person has to go back and insert all those changes in Markdown]
```

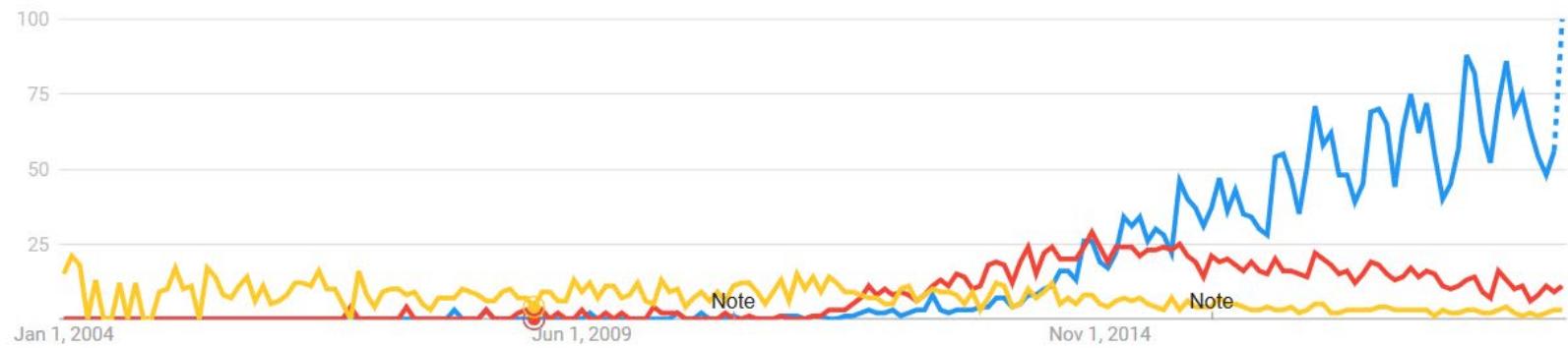
Right Pane (Rendered Output):

The right pane contains the rendered R Markdown content, which includes the title, author information, and the R code chunks. The rendered output shows the results of the R code, such as the summary of the cars dataset and the updated mean speed after changing one value.

Tools for Dynamic Documents and R

≡ Google Trends

Explore



● R Markdown
Search term

● knitR
Search term

● Sweave
Search term

Reproducible Research: Incidents raise awareness

The Annals of Applied Statistics
2009, Vol. 3, No. 4, 1309–1334
DOI: 10.1214/09-AOAS291
© Institute of Mathematical Statistics, 2009

DERIVING CHEMOSENSITIVITY FROM CELL LINES: FORENSIC BIOINFORMATICS AND REPRODUCIBLE RESEARCH IN HIGH-THROUGHPUT BIOLOGY

BY KEITH A. BAGGERLY¹ AND KEVIN R. COOPER²

This article has been retracted

An Expression of Concern has been published for this article

Science 12 December 2014:
Vol. 346 no. 6215 pp. 1366-1369
DOI: 10.1126/science.1256151

REPORT

When contact changes minds: An experiment for gay equality

Michael J. LaCour¹, Donald P. Green²

Author Affiliations

¹Department of Political Science, University of California, Los Angeles, CA, USA

²Department of Political Science, Columbia University, New York, NY, USA

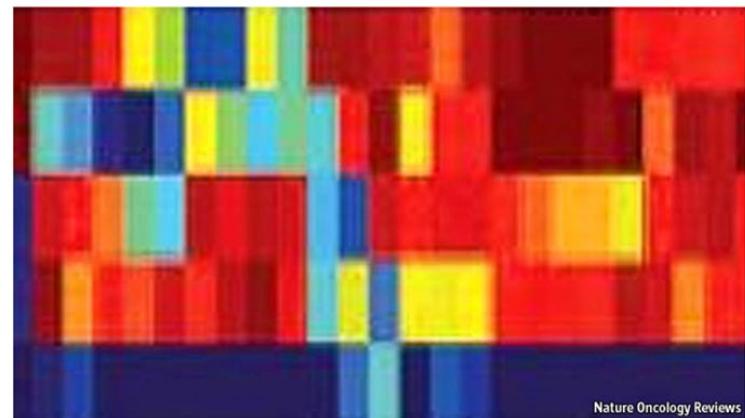
Misconduct in science

An array of errors

Investigations into a case of alleged scientific misconduct have revealed numerous holes in the oversight of science and scientific publishing

Sep 10th 2011 | From the print edition

 Like 962  Tweet 217



Nature Oncology Reviews

ANIL POTTI, Joseph Nevins and their colleagues at Duke University in Durham, North Carolina, garnered widespread attention in 2006. They reported in the *New England Journal of Medicine* that they could predict the course of a patient's lung cancer using devices called expression arrays, which log the activity patterns of thousands of genes in a sample of tissue as a colourful picture (see above). A few months later, they wrote in *Nature Medicine* that they had developed a similar technique which used gene expression in laboratory cultures of cancer cells, known as cell lines, to predict which chemotherapy would be most effective for an individual patient suffering from lung, breast or ovarian cancer.

At the time, this work looked like a tremendous advance for personalised medicine—the idea that understanding the molecular specifics of an individual's illness will lead to a tailored

Reproducible Research: NIH Statement

U.S. Department of Health & Human Services

NIH National Institutes of Health
Turning Discovery Into Health

Search NIH

NIH Employee Intranet | Staff Directory | En Español

Health Information Grants & Funding News & Events Research & Training Institutes at NIH About NIH

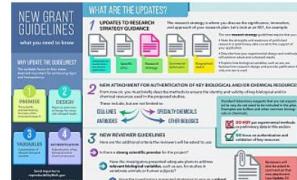
Home » Research & Training

RIGOR AND REPRODUCIBILITY



Rigor and Reproducibility

[Principles and Guidelines](#)
[Expanded Guidelines](#)
[Application Instructions](#)
[Training](#)
[Funding Opportunities](#)
[Meetings and Workshops](#)
[Announcements](#)
[Publications](#)



Two of the cornerstones of science advancement are rigor in designing and performing scientific research and the ability to reproduce biomedical research findings. The application of rigor ensures robust and unbiased experimental design, methodology, analysis, interpretation, and reporting of results. When a result can be reproduced by multiple scientists, it validates the original results and readiness to progress to the next phase of research. This is especially important for clinical trials in humans, which are built on studies that have demonstrated a particular effect or outcome.

In recent years, however, there has been a growing awareness of the need for rigorously designed published preclinical studies, to ensure that such studies can be reproduced. This webpage provides information about the efforts underway by NIH to enhance rigor and reproducibility in scientific research.



Email Updates

Sign up to receive email updates about rigor and reproducibility.

[Sign up for updates](#)

Related Links

[Letter from Dr. Stephen I. Katz: An Update on the NIH Initiative to Enhance Research Rigor and Reproducibility](#)

Contact Us

Please send email to NIHReprodEfforts@od.nih.gov.

Reproducible Research: Beyond Biostatistics

Reproducible vs Replicable

Using the **same (raw) data** and information about the analysis methods and choices, we can recreate the results.

Scientific finding is verified or supported by **independent experiment or study**.

“We define reproducibility as the ability to re-compute data analytic results given an observed dataset and knowledge of the data analysis pipeline. The replicability of a study is the chance that an independent experiment targeting the same scientific question will produce a consistent result.”

-- Leek and Peng 2015

Leek and Peng “Opinion: Reproducible Research can still be wrong: Adopting a prevention approach” PNAS, February 10, 2015, vol. 112, no. 6, 1645–1646

Reproducible Research: What can it really do?

Strengths and Limitations

Data Replication & Reproducibility

PERSPECTIVE

Reproducible Research in Computational Science

Roger D. Peng

Computational science has led to exposed limitations in our ability to serve as a minimum standard for study is not possible.

The rise of computational science exciting and fast-moving developments in many scientific areas. New techniques increased computing power, and methodological advances have dramatically improved to collect complex high-dimensional data sets. Large data sets have led to scientists computation, as well as researchers in traditionally oriented fields directly engage in science. The availability of large public datasets has allowed for researchers to make scientific contributions without using

Peng, R "Rep

computer. Making these computer codes available to others provides a level of detail regarding the analysis that is greater than the analogous non-computational experimental descriptions printed in journals using a natural language.

A critical barrier to reproducibility in many cases is that the computer code is no longer avail-

Opinion: Reproducible research can still be wrong: Adopting a prevention approach

Jeffrey T. Leek^{a,1} and Roger D. Peng^b

^aAssociate Professor of Biostatistics and Oncology and ^bAssociate Professor of Biostatistics, Johns Hopkins University, Baltimore, MD

Reproducibility—the ability to recompute results—and replicability—the chances other experimenters will achieve a consistent result—are two foundational characteristics of successful scientific research. Consistent findings from independent investigators are the primary means by which scientific evidence accumulates for or against a hypothesis. Yet, of late, there has been a crisis of confidence among researchers worried about the rate at which studies are either reproducible or replicable. To maintain the integrity of science research and the public's trust in science, the scientific community must ensure reproducibility and replicability by engaging in a more preventative approach that greatly expands data analysis education and routinely uses software tools.

There have been some very public failings of reproducibility across a range of disciplines from cancer genomics (3) to economics (4), and the data for many publications have not been made publicly available, raising doubts about the quality of data analyses. Popular press articles have raised questions about the reproducibility of all scientific research (5), and the US Congress has convened hearings focused on the transparency of scientific research (6). The result is that much of the scientific enterprise has been called into question, putting funding and hard won scientific truths at risk.

From a computational perspective, there are three major components to a reproducible and replicable study: (i) the raw data from the experiment are available, (ii) the statistical

computational tools such as knitr, iPython notebook, LONI, and Galaxy (8) have simplified the process of distributing reproducible data analyses.

Unfortunately, the mere reproducibility of computational results is insufficient to address the replication crisis because even a reproducible analysis can suffer from many problems—confounding from omitted variables, poor study design, missing data—that threaten the validity and useful interpretation of the results. Although improving the reproducibility of research may increase the rate at which flawed analyses are uncovered, as recent high-profile examples have demonstrated (4), it does not change the fact that problematic research is conducted in the first place.

The key question we want to answer when seeing the results of any scientific study is “Can I trust this data analysis?” If we think of problematic data analysis as a disease, reproducibility speeds diagnosis and treatment in

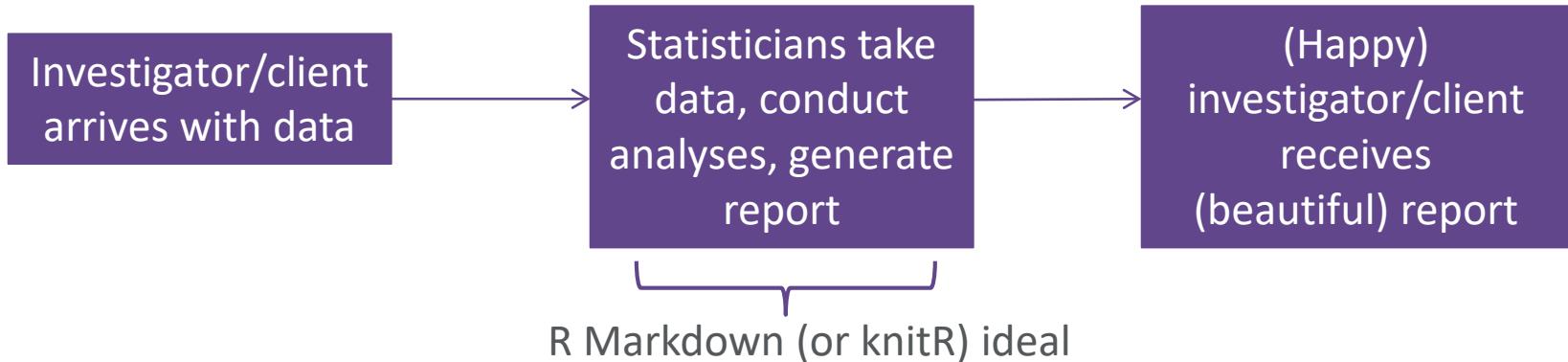
Leek and Peng “Opinion: Reproducible Research can still be wrong: Adopting a prevention approach” PNAS, February 10, 2015, vol. 112, no. 6, 1645–1646



Dynamic Documents: Challenges for Collaboration

Reproducible Research: Role of the Statistician

Challenges with Collaborations and Consultations



Reproducible +

Investigator/client comes back in a year with questions, we still know what we did

Investigator/client comes back with updated data, no problem to update

Reproducible -

Reproducibility only exists (to our knowledge) in the middle step

What if we have greater responsibilities and involvement?

MS Word is Ubiquitous for Manuscript Preparation

Clinical and Basic Sciences



"All text...should be in one double-spaced electronic document (preferably a Word Doc)"



"For submission and review, please submit the manuscript as a Word document. Do not submit your manuscript in PDF format."



"Science prefers to receive files in Word's .docx format."

R Markdown and Microsoft Word

.Rmd file can generate a .rtf (.docx) file

```
1 ---  
2 title: 'Association of Education with Anthropometr  
3 and Nutrition Examination Study 2013-2014'  
4 geometry: margin=.75in  
5 output: html_document  
6 sansfont: Calibri Light  
7 fontsize: 11pt  
8 ---  
9  
10 <style type="text/css">  
11 h1.title {  
12   font-size: 11pt;  
13   font-weight: bold;  
14 }  
15 </style>  
16  
17 ```{r setup, include=FALSE}  
18 knitr::opts_chunk$set(echo = FALSE)  
19 setwd("R:/NUCATS/NUCATS_Shared/BERDShared/Analysis")  
20 analysis<-read.csv("Analysis.csv")  
21 attach(analysis)  
22 library(tableone)  
23 library(knitr)  
24 options(digits=2)  
25  
26  
27 \usepackage[utf8]  
28  
29 **Introduction:** Education level has been shown to be associated with body mass index (BMI). Gender, race, marital status and metabolic characteristics may alter this association.  
30  
31 **Methods:** This study included adult ( $\geq 30$  years) participants from the 2013-2014 National Health and Nutrition Examination Study (NHANES). Education and demographic information were assessed by questionnaire and anthropometric measurements were taken by study personnel. Education was dichotomized based on matriculation to post-secondary education. Associations were estimated using T-tests or Wilcoxon rank sum tests for continuous data and chi-squared tests for categorical data. We examined the association between BMI and education level using multivariate linear regression.  
32  
33 ```{r, echo=FALSE}  
34 ps<- table(analysis$Postsecondary)  
35 percents<-100*prop.table(ps)  
36 model <- lm(BMXBMI ~ Postsecondary, data = analysis)  
37 BetaUnivariable <- summary(model)$coefficients[2,]  
38 LBUnivariable <- (summary(model)$coefficients[2, 1])  
39 UBUnivariable <- (summary(model)$coefficients[2, 1])  
40 myVars <- c("Gender", "Race", "RIDAGEYR", "Married", "Tableone")  
41 Tableone <- CreateTableone(data = analysis, vars=
```

Association of Education with Anthropometrics in US Adults: National Health and Nutrition Examination Study 2013-2014

Introduction: Education level has been shown to be associated with body mass index (BMI). Gender, race, marital status and metabolic characteristics may alter this association.

Methods: This study included adult (≥ 30 years) participants from the 2013-2014 National Health and Nutrition Examination Study (NHANES). Education and demographic information were assessed by questionnaire and anthropometric measurements were taken by study personnel. Education was dichotomized based on matriculation to post-secondary education. Associations were estimated using T-tests or Wilcoxon rank sum tests for continuous data and chi-squared tests for categorical data. We examined the association between BMI and education level using multivariate linear regression.

Results: Among 4808 participants, 2649 (55.1%) self-reported any post-secondary education. Post-secondary education was associated with lower BMI (Beta: -0.62, 95% CI: -1.03 to -0.21). After adjusting for gender, race, age, marital status, fasting glucose and total cholesterol, post-secondary education was no longer significantly associated with BMI (Beta: -0.19, 95% CI: -0.79 to 0.41).

Table 1. Association of Education with Participant Characteristics among 2013-2014 NHANES Participants

	No Postsecondary	Postsecondary	p
n	2159	2649	
Gender = Male (%)	1079 (50.0)	1208 (45.6)	0.003
Race (%)			<0.001
Mexican American	456 (21.1)	173 (6.5)	
Non-Hispanic Asia	161 (7.5)	411 (15.5)	

Practical Limitations

I send my collaborators the dynamic document, they send back:

Association of Education with-and Anthropometrics in US Adults: Results from the National Health and Nutrition Examination Study 2013-2014

Introduction: Education level ~~may be has been shown to be~~ associated with body mass index (BMI). Gender, race, marital status and metabolic characteristics may alter this association.

Methods: ~~We studied This study included~~ adult (≥ 30 years) participants from the 2013-2014 National Health and Nutrition Examination Study (NHANES). Education and demographic information were ~~self-reported assessed by questionnaire~~ and

Two bad choices:

1. Continue in Word, and lose reproducibility.
2. Re-enter all of their changes in my source “.rmd” file.

~~education had~~ Postsecondary education was associated with significantly lower BMI (Beta: -0.62, 95% CI: -1.03 to -0.21), ~~although~~. After adjusting for gender, race, age, marital status, fasting glucose and total cholesterol, ~~the association was no longer statistically significant~~ post-secondary education was no longer significantly associated with BMI (Beta: -0.19, 95% CI: -0.79 to 0.41).

Table 1. Association of Education with Participant Characteristics among $n = 4808$ 2013-2014 NHANES Participants

	No Postsecondary	Postsecondary	p
n	2159	2649	
Gender = Male (%)	1079 (50.0)	1208 (45.6)	0.003
Race (%)			<0.001
Mexican American	456 (21.1)	173 (6.5)	



Leah J Welty A few seconds ago
Leah, the spacing in this table is too far apart. Can you please make it single-spaced? And make the middle two columns wider?



Leah J Welty
Formatted: Font: Bold



Leah J Welty
Formatted: Font: Bold



Leah J Welty
Formatted: Font: Bold

Practical Limitations

Collaborators and plain text editors

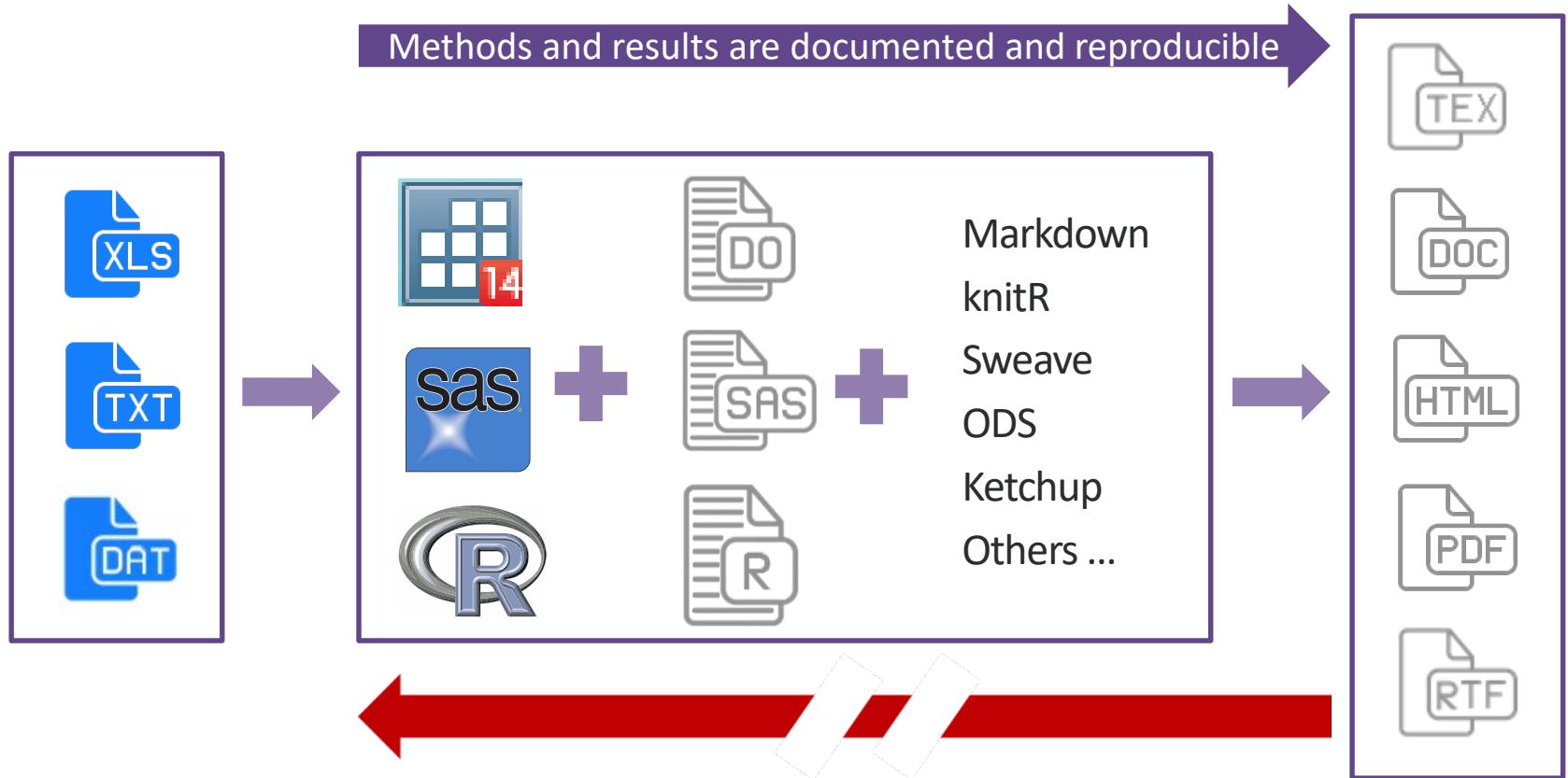
Markdown and knitR require drafting documents in a plain text editor:

```
3 author: "Leah Welty"
4 date: "July 27, 2006"
5 output: word_document
6 ---
7
8 ````{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10
11 You can use R Markdown from within RStudio. You write in a simple text editor, using the (fairly simple) Markdown language to indicate *italics* output in the document.
12
13 For example, if I want to see a summary of the *cars* dataset that comes standard with R, I can insert R code that produces this:
14
15 ````{r cars}
16 summary(cars)
17
18 I can also embed results directly in the text. For example, the median speed is `r mean(cars$speed)`.
19
20 That's pretty nice, because if I change something about the data, then that number can be automatically updated.
21
22 Another recent thing is that I can actually call and run Stata code from this interface. Neat, but I still have a problem ...
23
24
25
```

Are your collaborators willing to work this way?
My collaborators (primarily doctors and social scientists) are not.

Dynamic Documents: Existing Tools

Limitations of Existing Tools for Dynamic Documents



Why not R Markdown?

Markdown is GREAT, but not a universal solution for reproducibility

1. Collaborators may be unwilling to work in plain text.
2. Many (clinical) journals prefer/require submissions in MS Word.
3. Markdown can generate an RTF, but any changes to it must be re-entered in the source file.
4. Weaving tools (Markdown, Sweave, knitR) work with one software program at a time (e.g. R).
5. It can be messy to have statistical code and manuscript text all in the same place.



StatTag: Reproducible Research using Microsoft Word

Reproducible Research using Microsoft Word

Searching for tools in 2012

J Mesirov "Accessible Reproducible Research" *Science* 2010

Collaboration with Microsoft for Gene Pattern Software

Very limited use

Nolan, Peng and Lang 2011 "Enhanced Dynamic Documents for Reproducible Research"

RWordXML: cut and paste R code in to Word document

R package to process it and create a pdf or html file

Stata Automation Report

Word plug-in

Only for Stata

Worked only in some circumstances



The poster features a blue header bar with the text "VIII ITALIAN STATA USERS MEETING" and "Isola di San Servolo, Venezia November 17-18, 2011". Below the header is a small graphic of colored dots. In the center is a large gear icon. At the bottom, the text reads "Sar: Automatic Generation of Statistical Reports Using Stata and Microsoft Word for Windows". The footer contains the name "Giovanni Luca Lo Magno" with the email "lomagno.g@virgilio.it", the affiliation "Department of Economics, Business and Finance University of Palermo", and a note "Currently under review by the Stata Journal".

Finding a team and funding

Solving a pain-point takes resources



StatTag: A complementary approach

Incorporate existing strengths, pragmatic



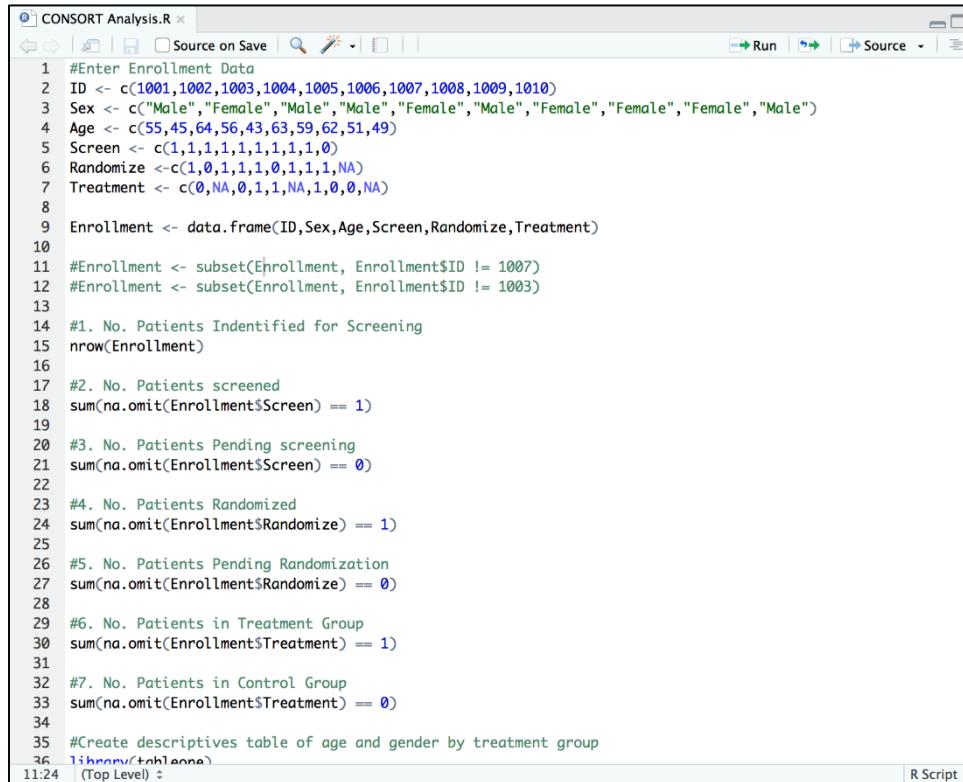
- StatTag creates a link between a statistical code file and a Word document
- Embeds output (values, tables, figures, verbatim) in document
- Can work separately on the code and the Word document but retain link
- Software agnostic: connects R, R Markdown (or Stata, SAS) to Word document
- Can connect multiple code files to the same document



StatTag: How it Works

StatTag

Step 1: Write your R/R Markdown code.



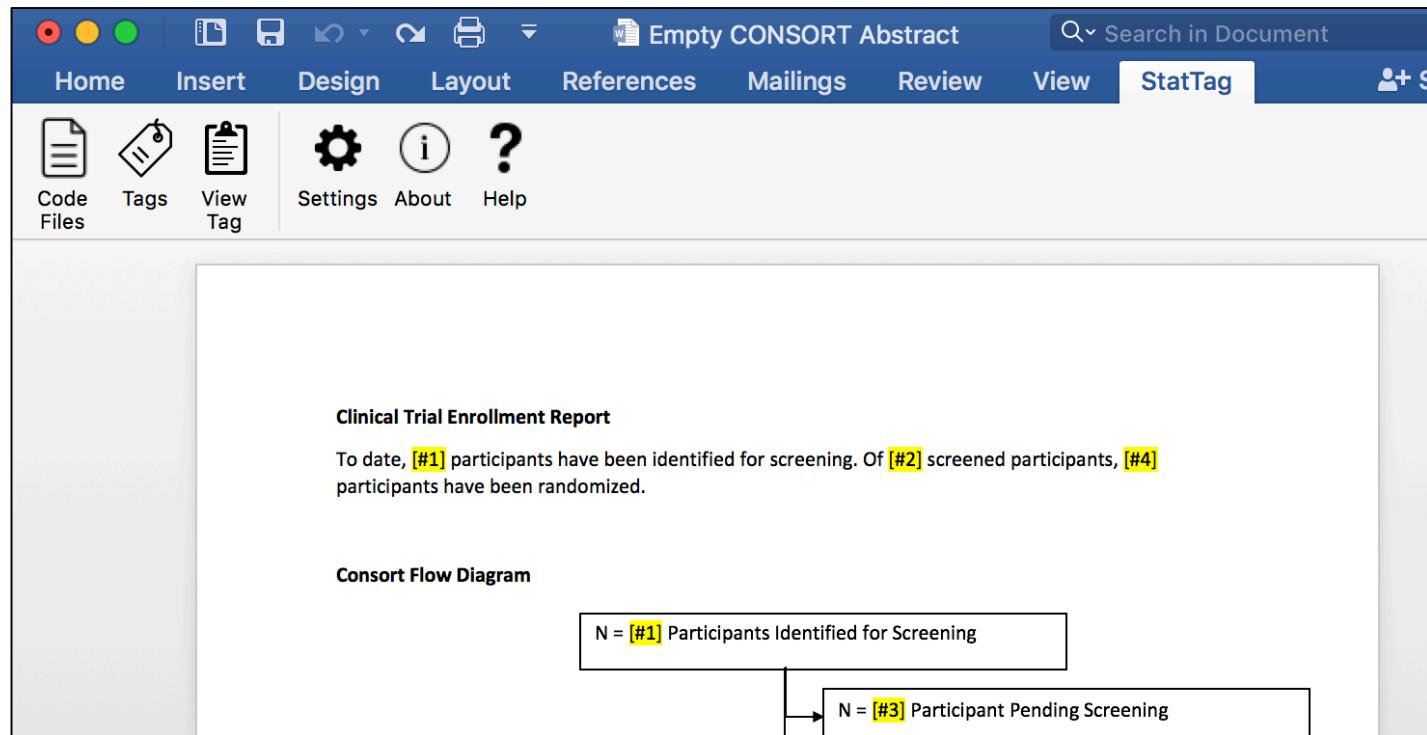
The screenshot shows an RStudio interface with the title bar "CONSORT Analysis.R". The code editor contains R code for performing a CONSORT analysis. The code is numbered from 1 to 36. It starts by defining variables for ID, Sex, Age, Screen, Randomize, and Treatment. It then subsets the Enrollment data frame to exclude rows where ID is 1007 or 1003. The code then counts the number of patients identified for screening (nrow(Enrollment)), the number of patients screened (sum(na.omit(Enrollment\$Screen) == 1)), the number of patients pending screening (sum(na.omit(Enrollment\$Screen) == 0)), the number of patients randomized (sum(na.omit(Enrollment\$Randomize) == 1)), the number of patients pending randomization (sum(na.omit(Enrollment\$Randomize) == 0)), the number of patients in the treatment group (sum(na.omit(Enrollment\$Treatment) == 1)), and the number of patients in the control group (sum(na.omit(Enrollment\$Treatment) == 0)). Finally, it creates a descriptive table of age and gender by treatment group using the knitr package.

```
1 #Enter Enrollment Data
2 ID <- c(1001,1002,1003,1004,1005,1006,1007,1008,1009,1010)
3 Sex <- c("Male","Female","Male","Female","Male","Female","Female","Male")
4 Age <- c(55,45,64,56,43,63,59,62,51,49)
5 Screen <- c(1,1,1,1,1,1,1,1,0)
6 Randomize <- c(1,0,1,1,1,0,1,1,1,NA)
7 Treatment <- c(0,NA,0,1,1,NA,1,0,0,NA)
8
9 Enrollment <- data.frame(ID,Sex,Age,Screen,Randomize,Treatment)
10
11 #Enrollment <- subset(Enrollment, Enrollment$ID != 1007)
12 #Enrollment <- subset(Enrollment, Enrollment$ID != 1003)
13
14 #1. No. Patients Identified for Screening
15 nrow(Enrollment)
16
17 #2. No. Patients screened
18 sum(na.omit(Enrollment$Screen) == 1)
19
20 #3. No. Patients Pending screening
21 sum(na.omit(Enrollment$Screen) == 0)
22
23 #4. No. Patients Randomized
24 sum(na.omit(Enrollment$Randomize) == 1)
25
26 #5. No. Patients Pending Randomization
27 sum(na.omit(Enrollment$Randomize) == 0)
28
29 #6. No. Patients in Treatment Group
30 sum(na.omit(Enrollment$Treatment) == 1)
31
32 #7. No. Patients in Control Group
33 sum(na.omit(Enrollment$Treatment) == 0)
34
35 #Create descriptives table of age and gender by treatment group
36 library(knitr)
37 knitr::tableone()
```

StatTag: How it Works

StatTag

Step 2: Open Word and write some or all of your text.



StatTag: How it Works

StatTag

Step 3: Use StatTag to connect your code file to the Word document.

```
CONSORT Analysis.R
Source on Save | Run | Source | R Script |
1 #Enter Enrollment Data
2 ID <- c(1001,1002,1003,1004,1005,1006,1007,1008,1009,1010)
3 Sex <- c("Male","Female","Male","Female","Male","Female","Female","Male")
4 Age <- c(55,45,64,56,43,63,59,62,51,49)
5 Screen <- c(1,1,1,1,1,1,1,1,1,0)
6 Randomize <-c(1,0,1,1,1,0,1,1,1,0)
7 Treatment <- c(0,NA,0,1,1,NA,1,0,0,NA)
8
9 Enrollment <- data.frame(ID,Sex,Age,Screen,Randomize,Treatment)
10
11 #Enrollment <- subset(Enrollment, Enrollment$ID != 1007)
12 #Enrollment <- subset(Enrollment, Enrollment$ID != 1003)
13
14 #1. No. Patients Identified for Screening
15 nrow(Enrollment)
16
17 #2. No. Patients screened
18 sum(is.na.omit(Enrollment$Screen) == 1)
19
20 #3. No. Patients Pending screening
21 sum(is.na.omit(Enrollment$Screen) == 0)
22
23 #4. No. Patients Randomized
24 sum(is.na.omit(Enrollment$Randomize) == 1)
25
26 #5. No. Patients Pending Randomization
27 sum(is.na.omit(Enrollment$Randomize) == 0)
28
29 #6. No. Patients in Treatment Group
30 sum(is.na.omit(Enrollment$Treatment) == 1)
31
32 #7. No. Patients in Control Group
33 sum(is.na.omit(Enrollment$Treatment) == 0)
34
35 #Create descriptive table of age and gender by treatment group
36 .11knewn/1hLambda
37 (Top Level) :
```

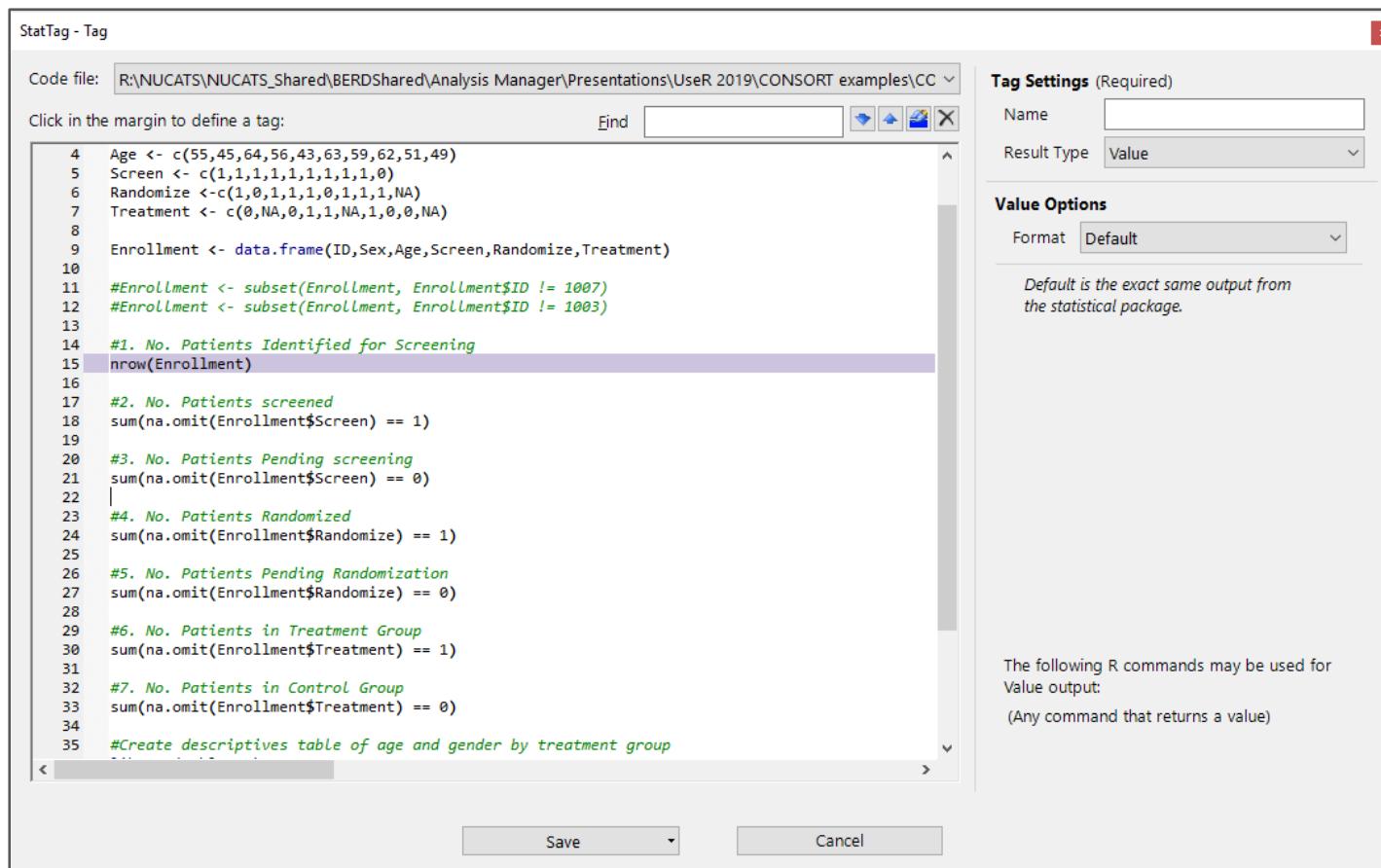
The screenshot shows a Microsoft Word document titled "Empty CONSORT Abstract". The ribbon at the top has a "StatTag" tab selected, which is highlighted with a purple circle. The "Home" tab is also visible. Below the ribbon, there are several icons: "Code Files" (document icon), "Tags" (tag icon), "View Tag" (document icon), "Settings" (gear icon), "About" (info icon), and "Help" (question mark icon). The main content area of the Word document contains a section titled "Clinical Trial Enrollment Report". It includes the following text:
To date, #[#1] participants have been identified for screening. Of #[#2] screened participants, #[#4] participants have been randomized.

Consent Flow Diagram
A flow diagram with two boxes connected by an arrow. The top box contains the text "N = [#1] Participants Identified for Screening". The bottom box contains the text "N = [#3] Participant Pending Screening". An arrow points from the top box to the bottom box.

StatTag: How it Works

StatTag

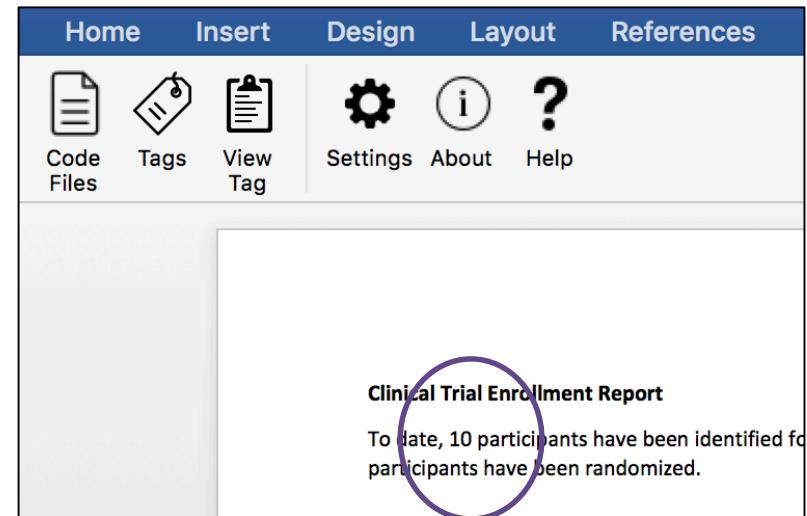
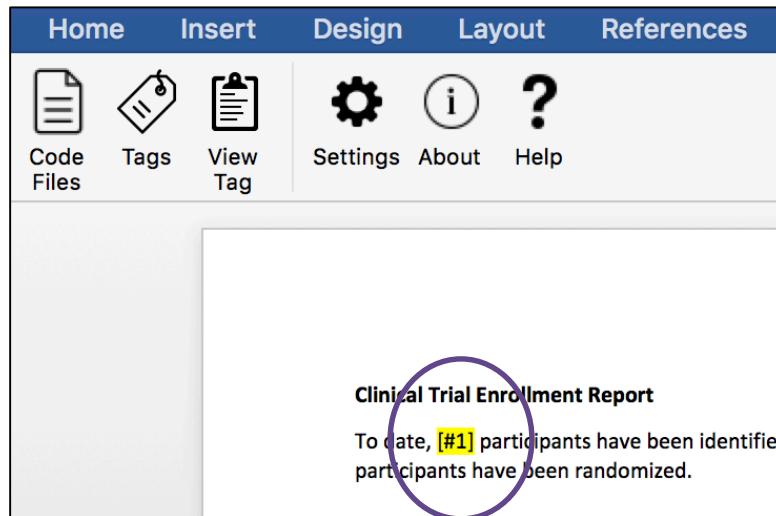
Step 4: Use StatTag to identify “tags” – portions of output that you want to insert – using a pop-up from Word.



StatTag: How it Works

StatTag

Step 5: Use StatTag to insert “tags” in your document. StatTag will run the code in the background and insert the results in to Word.



StatTag Demo video at Open StatTag on YouTube:

<https://www.youtube.com/watch?v=K3QwG4LB9a4>

StatTag: How it Works

Double-clicking a tag brings up associated code

StatTag

StatTag - Tag

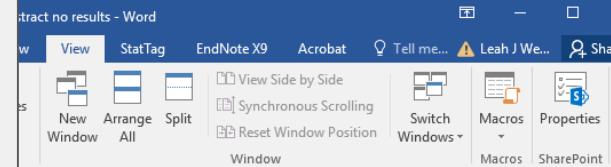
Code file: \\fsmresfiles\fsmresfiles\NUCATS\NUCATS_Shared\BERDShared\Analysis Manager\Presentations\Short Course

Click in the margin to define a tag:

```
28 #####ST:Value(Label="UnivariableBeta", Frequency="Always", Type="Numeric", Decimals=1, ThousandSeparator=",")
29 summary(model)$coefficients[2, 1]
30 ##<<<
31 #####ST:Value(Label="UniLB", Frequency="Always", Type="Numeric", Decimals=1, Thousands=False)
32 print((summary(model)$coefficients[2, 1]) - 1.96 * (summary(model)$coefficients[2, 2]))
33 ##<<<
34 #####ST:Value(Label="UniUB", Frequency="Always", Type="Numeric", Decimals=1, Thousands=False)
35 print((summary(model)$coefficients[2, 1]) + 1.96 * (summary(model)$coefficients[2, 2]))
36 ##<<<
37
38 ##Table 1
39 myVars <- c("Gender", "Race", "RIDAGEYR", "Married", "BMXBMI", "LBXTC", "LBXGLU")
40 TableOne <- CreateTableOne(data = analysis, vars=myVars, strata= "PostSecondary")
41 SkewedVars <- c("BMXBMI", "LBXTC", "LBXGLU")
42 print(TableOne, nonnormal = SkewedVars)
43
44 Table1 <- print(TableOne, nonnormal = SkewedVars, quote = FALSE, noSpaces = TRUE, printToggle=TRUE)
45
46 ## Save to a csv file
47 #####ST:Table(Label="Table 1", Frequency="On Demand", ColFilterEnabled=True, ColFilterType="Entire Column")
48 write.csv(Table1, file = "Table1.csv")
49 ##<<<
50
51 ## Figure 1.
52 #####ST:Figure(Label="Figure 1", Frequency="Always")
53 jpeg('figure1.jpg')
54 boxplot(BMXBMI~PostSecondary,data=analysis, main="BMI by Education Level",
55 ylab="Body Mass Index (kg/m**2)",names=c("Secondary", "Postsecondary"))
56 dev.off()
57 ##<<<
58
59 ## Final Model
```

Save Cancel

The following R commands may be used for Value output:
(Any command that returns a value)



US Adults: National Health and Nutrition Examination Study

associated with body mass index (BMI). Gender, race, and education level were significantly associated with this association.

Methods: This study included a total of 20 years) participants from the 2013-2014 National Health and Nutrition Examination Study (NHANES). Demographic and demographic information were assessed by questionnaire and anthropometric measurements were taken by study personnel. Education was dichotomized based on matriculation to post-secondary school. Associations were estimated using T-tests or Wilcoxon rank sum tests for continuous data and chi-square tests for categorical data. We examined the association between BMI and education level using multivariate regression.

Results: Among 4808 participants, 2649 (55%) self-reported any post-secondary education. Post-secondary education was associated with lower BMI (Beta: -0.6, 95% CI: -1.0 to -0.2). After adjusting for gender, race, age, marital status, fasting glucose and total cholesterol, post-secondary education was no longer significantly associated with BMI (Beta:-0.2, 95% CI: -0.8 to 0.4).

Table 1. Association of Education with Participant Characteristics among 2013-2014 NHANES Participants

Characteristic	No Post-Secondary n=2159	Post-Secondary n=2649	P-Value
Gender = Male (%)	1079 (50.0)	1208 (45.6)	<0.001
Race (%)			
Mexican American	456 (21.1)	173 (6.5)	
Non-Hispanic Asian	161 (7.5)	111 (4.5)	
Non-Hispanic Black	102 (4.7)	408 (15.3)	
Non-Hispanic White	640 (29.4)	537 (20.1)	

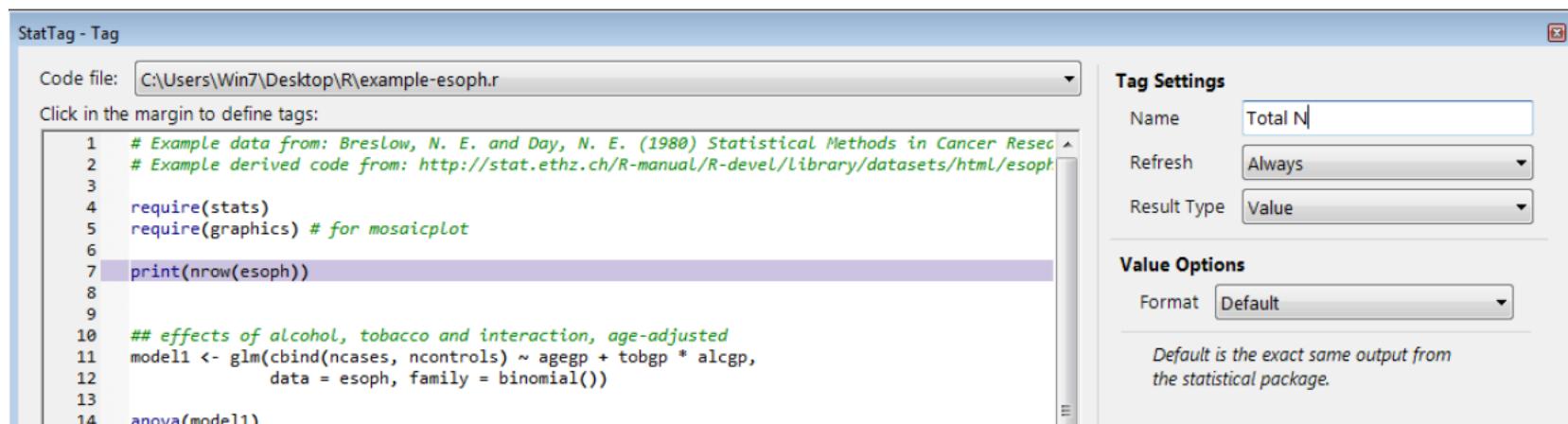
StatTag: How it works

Inserting tags directly in code, or using StatTag dialog box

StatTag

```
**>>>ST:Value(Label=" ", Frequency="", Type="")
[R, SAS, or Stata code]
**<<<

**>>>ST:Table(Label="", Frequency="", Type="", AllowInvalid=True,
Decimals=0, Thousands=False)
[R, SAS, or Stata code]
**<<<
```



StatTag: How it Works

StatTag

Recognizes different key words to identify results.

Type			or
Numeric Values	display	%put	Any command that returns a value
Tables	matrix list, any command that returns a .xls or .csv file	ODS CSV	Any command that returns a data frame, matrix, vector or list, any command that returns a .xls or .csv file
Figures	graph export	ODS PDF	pdf, win.metafile, png, jpeg, bmp, postscript
Verbatim	Any Code	Any Code	Any Code

StatTag: How it Works

Sharing Word documents (and code) with collaborators

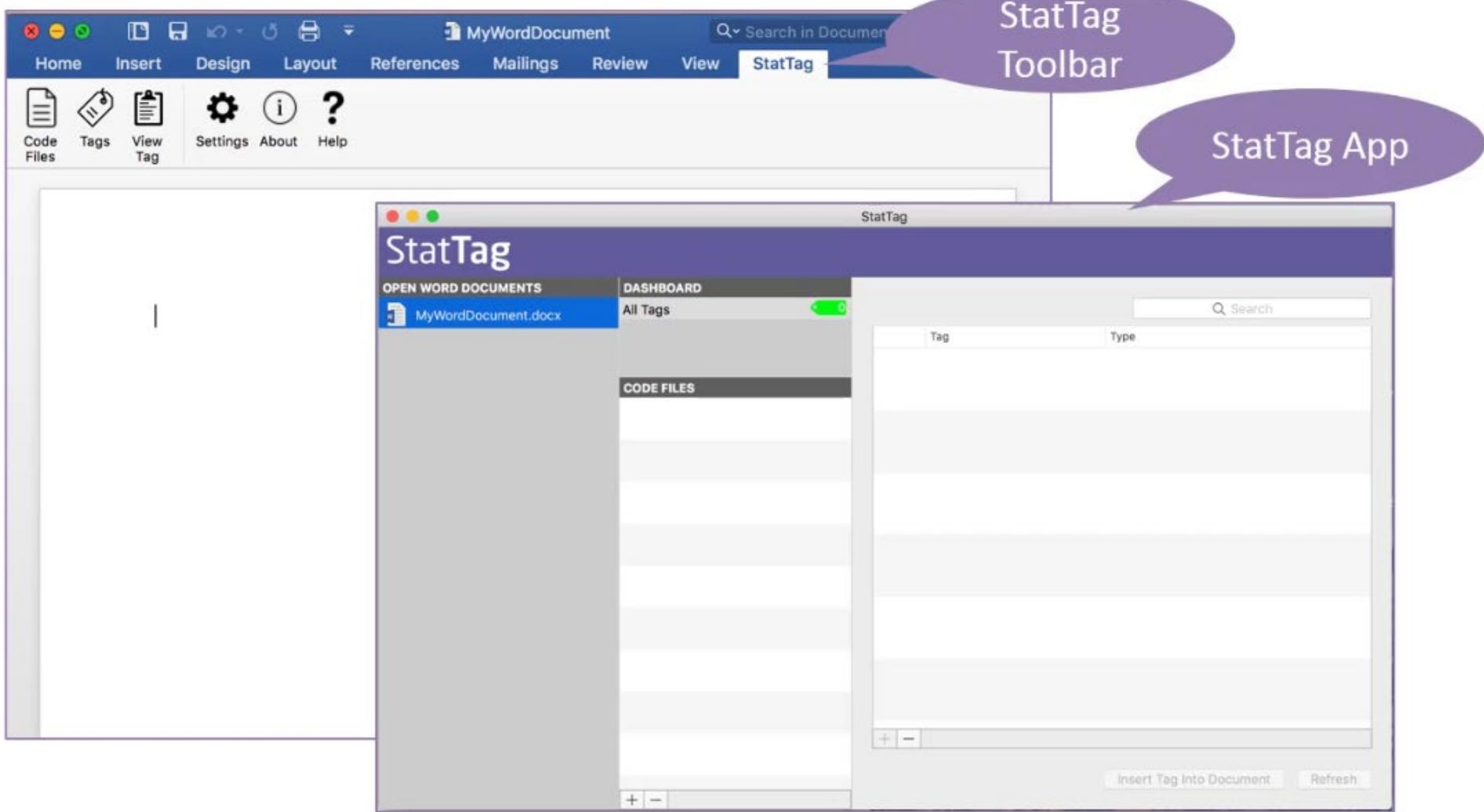
StatTag

If I have...	I can...		
	Review/edit manuscript text	View code associated with a tag	Insert or update a tag
Microsoft Word	✓	✗	✗
+ StatTag and Stata/SAS/R code	✓	✓	✗
+ Stata/SAS/R code and Data	✓	✓	✓

StatTag: How it Works (Mac)

App, rather than plug-in

StatTag



Getting StatTag

Freely available, open source, and evolving

StatTag

stattag.org

StatTag

Northwestern
University

[download stattag](#) / [user guide and tutorial](#) / [cite stattag](#) / [announcements](#) / [faq](#) / [contact](#)

STATTAG

StatTag is a free software plug-in for conducting reproducible research. It facilitates the creation of dynamic documents using Microsoft Word documents and statistical software, such as Stata. Users can use StatTag to embed statistical output (estimates, tables and figures) into a Word document and then with one click individually or collectively update output with a call to the statistical program. What makes StatTag different from other tools for creating dynamic documents is that it allows for statistical code to be edited directly from Microsoft Word. Using StatTag means that modifications to a dataset or analysis no longer require transcribing or re-copying results into a manuscript or table.

(Really) Getting StatTag

GitHub repo at github.com/StatTag



Screenshot of the GitHub organization page for StatTag.

Header: Features, Business, Explore, Marketplace, Pricing, This organization, Search, Sign in or Sign up.

Profile Section: StatTag, Chicago, IL, <http://stattag.org>.

Repository Buttons: Repositories (selected), People (0).

Search and Filters: Search repositories..., Type: All, Language: All.

Simple-Code-Examples Repository: Easy code examples for StatTag, updated on May 18 by SAS.

StatTag Repository: Windows version of StatTag, updated on May 18 by SAS.

Metrics: Top languages: C# (green), SAS (red), C++ (pink). People: 0. This organization has no public members. You must be a member to see who's a part of this organization.

Learning StatTag

Open course on Instructure Canvas

StatTag

The screenshot shows the StatTag 101 course page on the Canvas platform. The left sidebar contains a navigation menu with options like Home, Modules, Quizzes, Assignments, Files, People, Syllabus, Outcomes, Conferences, Grades, Pages, Collaborations, Discussions, Announcements, and Settings. The 'Home' button is highlighted in blue. The main content area displays the course title 'StatTag 101' and a 'Welcome!' message: 'We are excited you are interested in reproducible research. This course comprises exercises that will introduce you to StatTag.' Below this is an 'Instructions' section with a note about the online course content being arranged in Modules. It recommends starting with the 'Prework' module, which contains instruction to download and install StatTag as well as materials for the exercises. A table lists the modules and their content:

Module	Content
About the Course	General overview of the course
Prework	Prework exercises to be completed before the course
CONSORT Abstract	Hands on exercises for the course

On the right side, there is a 'Course Status' section showing 'Published' (with a green checkmark) and several other buttons for managing the course: Import from Commons, Choose Home Page, View Course Stream, Course Setup Checklist, New Announcement, Student View, and View Course Analytics. Below these are sections for 'Coming Up' (nothing for the next week) and 'Recent Feedback' (nothing for now).

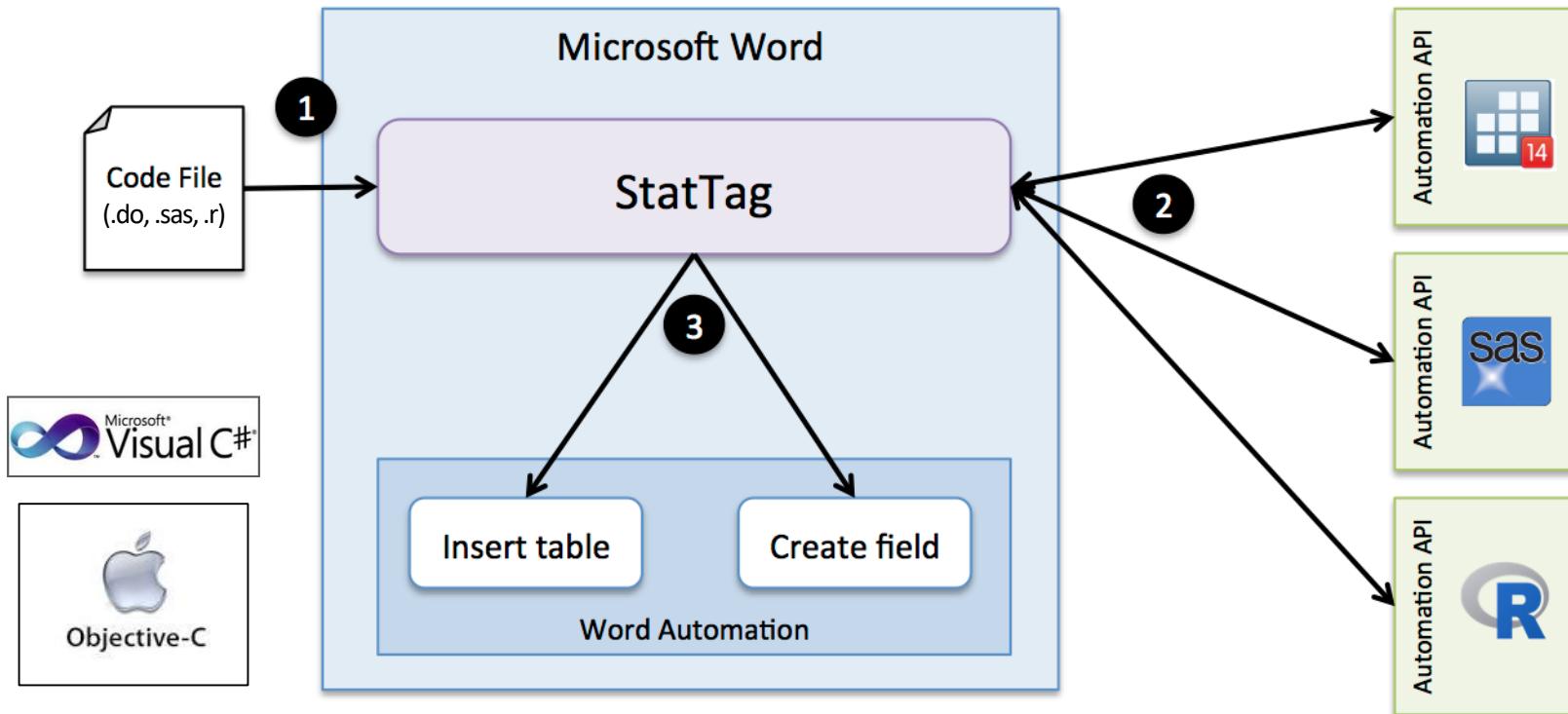
Online course available at: <https://canvas.instructure.com/enroll/DPCLGC>

Or enter code DPCLGC at <https://canvas.instructure.com/register>

StatTag: Architecture

How we built it

StatTag



1. Read the code file & parse out the tags
2. Send commands to the stat program and get individual results
3. Use Word automation to add results to the document (using native Word formatting for tables and fields).

Future directions

Tech savvy tools to adopt for collaborative reproducibility

- Connecting to even more
 - Python
 - Jupyter notebook



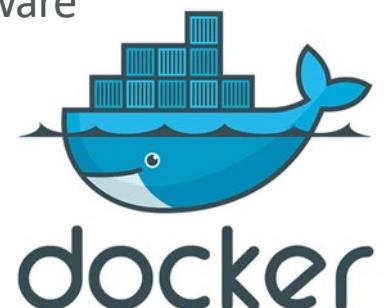
IP[y]:
IPython

- Version Control
 - Github
 - Subversion, cvs, etc.
 - Great tools, but what about WYSISYG or non-plain text files?



SUBVERSION®

- “Containers” (e.g. docker)
 - Facilitate packaging the code, data, packages, libraries and software program/environment together



Citations and Acknowledgements

Thank you!

- Welty, L.J., Rasmussen, L.V., Baldridge, A.S, and Whitley E. (2016). *StatTag*. Chicago, Illinois, United States: Galter Health Sciences Library.
doi:10.18131/G3K76
- StatTag is being developed with funding through a Clinical Translational Sciences Award (CTSA) to Northwestern University.





Questions?

- Visit stattag.org for more information
- Get involved! Contact us at stattag@northwestern.edu
 - Questions
 - Comments and suggestions
 - Volunteer to test new releases – *especially if you are also a Python user!*