

# Modern R in a Corporate Environment

Original materials developed for RADARS

*Brian Davis*

*2018-04-18*



# Contents

<b>About</b>	<b>5</b>
<b>1 Introduction</b>	<b>7</b>
1.1 Course Philosophy . . . . .	7
1.2 Prerequisites . . . . .	8
1.3 Content . . . . .	8
1.4 Structure . . . . .	9
<b>2 R Programming Basics</b>	<b>11</b>
2.1 Objects . . . . .	11
2.2 Important operators and assignment . . . . .	11
2.3 Comparison . . . . .	11
2.4 Basic math . . . . .	11
2.5 Logical and Sets . . . . .	11
2.6 Control flow . . . . .	11
2.7 Ordering and tabulating . . . . .	11



# About



# Chapter 1

## Introduction

Something that will make life easier in the long-run can be the most difficult thing to do today. For coders, prioritising the long term may involve an overhaul of current practice and the learning of a new skill.

### 1.1 Course Philosophy

“The best programs are written so that computing machines can perform them quickly and so that human beings can understand them clearly. A programmer is ideally an essayist who works with traditional aesthetic and literary forms as well as mathematical concepts, to communicate the way that an algorithm works and to convince a reader that the results will be correct.” Donald Knuth

#### 1.1.1 Reproducible Research

Reproducible research is the idea that data analyses, and more generally, scientific claims, are published with their data and software code so that others may verify the findings and build upon them. There are two basic reasons to be concerned about making your research reproducible. The first is *to show evidence of the correctness of your results*. The second reason to aspire to reproducibility is *to enable others to make use of our methods and results*.

Modern challenges of reproducibility in research, particularly computational reproducibility, have produced a lot of discussion in papers, blogs and videos, some of which are listed [here](#) and [here](#).

Conclusions in experimental psychology often are the result of null hypothesis significance testing. Unfortunately, there is evidence ((from eight major psychology journals published between 1985 and 2013) that roughly half of all published empirical psychology articles contain at least one inconsistent p-value, and around one in eight articles contain a grossly inconsistent p-value that makes a non-significant result seem significant, or vice versa. [statscheck](#) and [here](#)

“A key component of scientific communication is sufficient information for other researchers in the field to reproduce published findings. For computational and data-enabled research, this has often been interpreted to mean making available the raw data from which results were generated, the computer code that generated the findings, and any additional information needed such as workflows and input parameters. Many journals are revising author guidelines to include data and code availability. We chose a random sample of 204 scientific papers published in the journal **Science** after the implementation of their policy in February 2011. We found that were able to

reproduce the findings for 26%.” Proceedings of the National Academy of Sciences of the United States of America

“Starting September 1 2016, JASA ACS will require code and data as a minimum standard for reproducibility of statistical scientific research.” JASA

### 1.1.2 FDA Validation

“Establishing documented evidence which provides a high degree of assurance that a specific process will consistently produce a product meeting its predetermined specifications and quality attributes.” -Validation as defined by the FDA in **Validation of Systems for 21 CFR Part 11 Compliance**

### 1.1.3 The SAS Myth

Contrary to what we hear the FDA does not require SAS to be used *EVER*. There are instances that you have to deliver data in XPORT format though which is open and implemented in many programming languages.

“FDA does not require use of any specific software for statistical analyses, and statistical software is not explicitly discussed in Title 21 of the Code of Federal Regulations [e.g., in 21CFR part 11]. However, the software package(s) used for statistical analyses should be fully documented in the submission, including version and build identification. As noted in the FDA guidance, E9 Statistical Principles for Clinical Trials” FDA Statistical Software Clarifying Statement

Good write up with links to several FDA talks on the subject.

## 1.2 Prerequisites

- We will assume you have minimal experience and knowledge of R
- IT should have installed:
  - R version 3.5
  - RStudio version 1.1
  - MiTeX
  - RTools version 3.4
- We will install other dependencies throughout the course.

## 1.3 Content

It is impossible to become an expert in R in only one course even a multi-week one. Yet, this course aims at giving a wide understanding on many aspects of R as used in a corporate / production environment. It will roughly be based on R for Data Science. While this is an *excellent* resource it does not cover much of what we will need on a routine basis. Some external resources will be referred to in this book for you to be able to deepen what you would have learned in this course.

This is your course so if you feel we need to hit an area deeper, or add content based on a current need, let me know and we will work to adjust it.

The **rough** topic list of the course:

1. Good programming practices
2. Basics of R Programming
3. Importing Data



4. Tidying Data
5. Visualizing Data
6. Functions
7. Strings
8. Dates and Time
9. Communicating Results

Making Code Production Ready:

10. Functions (part II)
11. Assertions
12. Unit tests
13. Documentation
14. Communicating Results (part II)

## 1.4 Structure

My current thoughts are to meet an hour a week and discuss a topic. We will not be going strictly through the R4DS, but will use it as our foundation into the topic at hand. Then give an assignment due for the next week which we go over the solutions. We will incorporate these assignments into a RADARS R package(s?) so we will have a collection of usefull reusable code for the future.

Open to other ideas as we go along. I'm going to try to keep the assignments related to our current work (maybe working through Site Investigator and/or Subscriber Reports) so we can work on the class during work hours.



## Chapter 2

# R Programming Basics

See this [vocabulary list](#) for a good starting point of the basics of base R and some important libraries.

### 2.1 Objects

#### 2.1.1 Making vectors

#### 2.1.2 Lists & data.frames

#### 2.1.3 Vectors and matrices

automatic coercion rules `character > numeric > logical`

### 2.2 Important operators and assignment

### 2.3 Comparison

### 2.4 Basic math

### 2.5 Logical and Sets

### 2.6 Control flow

### 2.7 Ordering and tabulating