

Newman Library





Weekly Materials Week Three - R Character Manipulation and Date Processing [02/13 - 02/19]

# Week Three - R Character Manipulation and Date Processing [02/13 - 02/19]



#### Week 3 Overview

This week, we'll work with character strings and dates in R.

#### **Events and Deliverables**

Thursday 02/16	Meetup, 6:30 p.m. EST
Thursday 02/16	Discussion Post due
Sunday 02/19	Week 3 Assignment due
Sunday 02/19	Discussion Response due
Sunday 09/26	Project 1 due



#### Week 3 Reading

Please read:

- Automated Data Collection in R, chapter 8.
- Data Science for Business, chapter 3.

Here is are some additional resources that you may find helpful.

- R for Data Science, chapters 11 ("Strings with stringr") and 13 ("Dates and Times with lubridate").
- "Handling and Processing Strings in R," Gaston Sanchez, http://gastonsanchez.com/Handling and Processing Strings in R.pdf, Includes link to freely downloadable (113 pp.!) eBook with good material on base R and stringr package string manipulation and grep.



#### **Regular Expressions**

Regular expressions provide advanced text processing capabilities. Often, you can choose between using R's string manipulation functions and "regex" functions. For messy text based source data, like scraping text from HTML-based web pages, regular expressions are often the best way forward.

Regular expressions are implemented somewhat differently in different programming languages, so you'll need to carefully test any regex code that you bring from another environment into R.

You may want to consult the overview article, "Regular Expressions as used in R", <a href="https://stat.ethz.ch/R-manual/R-devel/library/base/html/regex.html">https://stat.ethz.ch/R-manual/R-devel/library/base/html/regex.html</a>

The best introduction that I've found to Regular Expressions was put together by Software Carpentry for its (open source) Python course. Please also watch the video below:

Source: Software Carpentry, Inc. <a href="https://www.youtube.com/watch?v=c-0v1JUMDv4">https://www.youtube.com/watch?v=c-0v1JUMDv4</a>



Optional. If you're interested in learning more about Regular Expressions, here are links to the other Software Carpentry videos on regular expressions, with example code in Python.

- Regular Expresssion 2: Operators, <a href="https://www.youtube.com/watch?">https://www.youtube.com/watch?</a>
   y=G7 HnivvnyE
- Regular Expressions 3: Mechanics, <a href="https://www.youtube.com/watch?v=iixnLh55wpo">https://www.youtube.com/watch?v=iixnLh55wpo</a>
- Regular Expressions 4: Patterns, <a href="https://www.youtube.com/watch?v=FgxQyukp39A">https://www.youtube.com/watch?v=FgxQyukp39A</a>
- Regular Expressions 5: More Tools, <a href="https://www.youtube.com/watch?v=RGN5tS-2Zmo">https://www.youtube.com/watch?v=RGN5tS-2Zmo</a>



### Discussion 3: Analyzing sets of text/log files

For this week's data science context discussion, you're asked to identify (or imagine) a "use case" that involves "loading a set of text files, and then performing

an analysis on the loaded, combined dataset."

A short description of your use case is sufficient. Here is an example:

"Each day, my organization's primary database creates a log file with about 50 MB of activity information. I would create an automated process to load the log files generated by my database server, and generate a daily report that tells me what are the five most frequently run queries and five slowest running queries for the previous day, and I would track how these frequent and/or slow queries change over time. This would help me to prioritize my database performance tuning efforts, such as deciding when to add an index to a specific table."

Please post by end of day on Thursday February 16<sup>th</sup>, and weigh in on at least one class mate's post by end of day on Sunday February 19<sup>th</sup>.

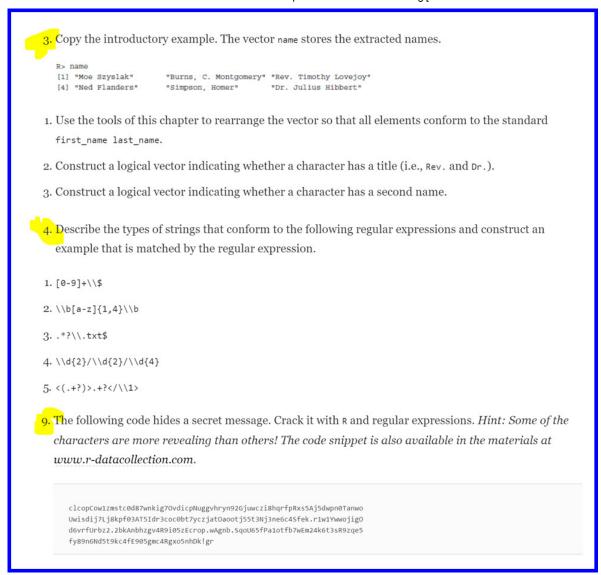


#### Week 3 Assignment

Please deliver links to an R Markdown file (in GitHub and rpubs.com) with solutions to problems 3 and 4 from chapter 8 of *Automated Data Collection in R*. Problem 9 is extra credit. You may work in a small group, but please submit separately with names of all group participants in your submission.

Here is the referenced code for the introductory example in #3:

raw.data <-"555-1239Moe Szyslak(636) 555-0113Burns, C. Montgomery555-65 Due end of day Sunday February  $19^{th}$ .





## Project 1

Attached Files: Project 1.pdf (120.332 KB)
tournamentinfo.txt (17.418 KB)

This project is due on Sunday February 26<sup>th</sup>. We'll review the most interesting solutions in our meetup on Thursday March 9<sup>th</sup>. (There is no meetup on Thursday March 2<sup>nd</sup>]

As before, please deliver your code in GitHub, with submission links to your GitHub repository and your published R Markdown file. You may work in a small group on this project.

Here is a 3 minute video on reading chess tournament cross-tables that you may find helpful:

# Reading a Chess Tournament Cross Table



[Chess ELO ratings have been adopted for use in ranking players and teams in other sports, as well as in a few human resources applications. You can learn more about ELO rating use cases here: <a href="https://fivethirtyeight.com/tag/elo-ratings/">https://fivethirtyeight.com/tag/elo-ratings/</a>]