



Data Management Techniques, Regular Expressions, and Dates and Times

James Balamuta

Department of Informatics, Statistics
University of Illinois at Urbana-Champaign

July 10, 2017

CC BY-NC-SA 4.0, 2016 - 2017, James J Balamuta

On the Agenda

1 Data Management Techniques

- Project Structure
- Reading and Writing
- Tidy Data (Reshaping)

2 character/string manipulation in R

- Counting, extracting, and concatenating
- Formatting with `format` and `sprintf`

3 Regular Expressions (regex)

- Operators
- How to extract information from character/strings?

4 Dates and Times

- Different Formats
- Converting to a `POSIXct` object
- `anytime` and `lubridate`

On the Agenda

1 Data Management

- Motivation
- Sample Project
- Case Study: Housing Data
- Tidy Data

2 Regular Expressions

- Motivation
- Syntax
- Case Study: House Addresses
- Extracting and Formatting

3 Dates and Times

- Motivation
- System Information
- Operations on Time
- Date Formats
- Time Formats

4 Misc

- POSIXlt
- anytime
- lubridate

Data Management Techniques

- When talking about data management techniques, we're going to aim to make things reproducible as always.
- Thus, everything should be done within a scripted context.
- To advocate this approach, we're going to be looking at a sample project setup and cleaning routine.

Example Project Setup & Files

- To follow along with the lecture, please visit the lecture 11 assets section of the class website to obtain the example project
- **Link**
 - <https://github.com/stat385uiuc/su2017/tree/master/static/assets/lectures/lec11/project>

Ideal Project Setup

Consider the following project directory structure:

```
| - Project
  | - data-raw/
    | - data.csv
    | - clean_data.R
  | - data/
    | - cleaned_data.rda
  | - R/
    | - analysis_script.R
  | - README.Rmd
  | - Project.Rproj
```

Notes:

- Anything related to cleaning data is within the data-raw directory.
- R code is contained to its own folder.
- README.Rmd file helps with understanding project contents.

OS Independent Load

Often times, we collaborate with a colleague through remote storage options such as [Dropbox](#) and [BoxSync](#). However, if your colleague is running a Mac and you are running Windows, how can you keep the same source files?

```
# Example of an independent OS environment
os_name = Sys.info()[['sysname']]
if(os_name == "Windows") {
  fp = file.path("F:/BoxSync/stat385/lectures")
} else if(os_name == "Darwin") { # macOS
  fp = file.path("~/BoxSync/stat385/lectures")
} else { # Linux
  stop("I'm a penguin.")
}

script_path = file.path(fp, "lec11")
```

Package Dependencies

Similarly, when sharing code, how can you make sure that the receiving party is able to run it, without having to look through the code to figure out dependencies (i.e. `library()` or `require()`)?

```
# Any package that is required by the script  
# below is given here  
inst_pkgs = load_pkgs = c("MASS", "faraway", "Hmsic",  
                          "randomForest", "rpart")  
  
# Check to see if the packages are already installed  
inst_pkgs = inst_pkgs[!(inst_pkgs %in%  
                        installed.packages()[, "Package"])]  
  
# Installs any missing package  
if(length(inst_pkgs)) install.packages(inst_pkgs)  
  
# Dynamically load required packages  
pkgs_loaded = lapply(load_pkgs, require, character.only=T)
```


Reading Data - R Core

The default ways to ingest data into *R* are:

- **base R:** Text Files

- `read.table()`: Read file in table format.
- `read.csv()` (US) / `read.csv2()` (Euro): Reads csv file
- `read.delim()` (US) / `read.delim2()` (Euro): Reads delimited file
- **Note:** Decimal separators differ depending on locale
 - e.g. USA uses `,` whereas Europe uses `.` hence `csv/csv2`

- **foreign:** Minitab, S, SAS, SPSS, Stata, Systat, Weka, dBase, Octave

Reading Data - Third Party

Two years ago, Hadley took it upon himself to write a lot of different file readers. In turn, these readers should be preferred as they are: 1. Faster and 2. Reliable

- **readr:** Text (.csv)
 - `read_file()`: Handles zipped files and .txt
 - `read_csv()`: CSV
- **readxl:** Excel
 - `read_excel()`: Reads in the first sheet (can be specified to others by name)
- **haven:** SAS, SPSS, and Stata
 - `read_sas()`: Handles SAS .b7dat and .b7cat
 - `read_spss()`: Handles SPSS .por and .sva
 - `read_stata()/read_dta()`: Handles Stata .dta

Load Multiple Data Sets

- Rarely is data ever in one source. Sometimes, data is found in a combination of **.csv** files with different filenames.

```
# Obtain a list of all files within active directory  
# with extension .csv  
filenames = list.files(pattern="*.csv")  
  
dsAll = data.frame()  
for (i in seq_along(filenames)) {  
  # Assign each file to its filename  
  assign(filenames[i],  
        read.csv(filenames[i], stringsAsFactors=FALSE))  
  # Note: The stringsAsFactors = FALSE condition  
  # prevents r from building levels into each variable.  
  
  # Quick bind (bad implementation)  
  dsAll = rbind(dsAll, filenames[i])  
}
```

Sample Data Clean

- For the next section, we will be focusing on the cleaning and formatting of data.
- Many of the ideas presented next were discussed by my good friend **Michael Quinn**.
- Michael obtained a Masters in Statistics from UIUC while working as a MAGNET Intern at the State Farm RDC and now works on Google's Ad Team!

Make a Rejection List

- When cleaning data, you will want to exclude an observation given a variable's specific value.
- To handle such values, we create a '**Rejection List**.' The list will specify observations that should be removed from the data.
 - Here is an example of a rejection list called `reject_list.txt`:

```
1002 w springfield ave  
315 s state st  
508 e soughton st  
....
```

Updating an Observation

- To update an observation, we look for a unique key to the observation.
 - In this case, it would be the house address. In other cases, it may be the Subject ID that is assigned.
- For most of the cases, we will want to only update a part of an observation.

```
# Find the row value
k = which(houses$St_Address == tolower("3909 Aberdeen Dr"))
# Update row value traits
houses$Bedrooms[k] = "3 beds"
houses$Bathrooms[k] = "2 baths"
houses$lot[k] = "3040 sqft"
houses$lastsoldifavailable[k] = "May 2011 for $135,000"
```

Efficient Vectorized Cleaning

- To clean efficiently means to vectorize your cleaning approaches:

```
# Removes any of the St_Addresses in the rejection list
houses = houses[!(houses$St_Address %in% reject),]
# Removes any duplicates
houses = unique(houses)
# Drop column range starting at heatttype to floorcover
# heatttype, zillowdays, cooling, parking, basement,
# fireplace, floorcover.
start_loc = match("heatttype",names(houses))
end_loc = match("floorcover",names(houses))
houses = houses[,-(start_loc:end_loc)]
# Removes variable headers in dataset
houses = houses[-which(houses$lastsoldifavailable
                        == "lastsoldifavailable"),]
```

Concept of Tidy Data

Tidy datasets are all alike but every messy dataset is messy in its own way

— Hadley Wickham (*JSS Tidy data*)

In tidy data:

- 1 Each variable forms a column.
- 2 Each observation forms a row.
- 3 Each type of observational unit forms a table.

Data Shape

Consider the following data set:

```
experiment = read.table(header=TRUE, text='
  subject sex control a b
      S1   F      4.2 4.1 2.2
      S2   M      5.9 7.2 6.8
      S3   M      9.1 9.8 10.2
      S5   F      2.1 23.5 5.2
')
```

experiment

```
##   subject sex control    a    b
## 1      S1   F      4.2 4.1 2.2
## 2      S2   M      5.9 7.2 6.8
## 3      S3   M      9.1 9.8 10.2
```

Wide Data

```
experiment
```

##	subject	sex	control	a	b
## 1	S1	F	4.2	4.1	2.2
## 2	S2	M	5.9	7.2	6.8
## 3	S3	M	9.1	9.8	10.2
## 4	S5	F	2.1	23.5	5.2

- In its current form, the data is considered to be **wide**.
- **Wide Data** has repeated responses or treatments of a subject in a single row with each response in its own column along with its properties.

Long Data

##	subject	sex	condition	measurement
## 1	S1	F	control	4.2
## 2	S2	M	control	5.9
## 3	S3	M	control	9.1
## 4	S5	F	control	2.1
## 5	S1	F	a	4.1
## 6	S2	M	a	7.2
## 7	S3	M	a	9.8
## 8	S5	F	a	23.5
## 9	S1	F	b	2.2
## 10	S2	M	b	6.8
## 11	S3	M	b	10.2
## 12	S5	F	b	5.2

- With a little modification, the data is considered to be **long**.
- **Long Data** has each row as one response per subject and any variables for the subject that do not change over time or treatment will have the same value in all the rows.

Wide Data to Long

- Use gather to move to a long format

```
library(tidyr)
(data_long = gather(experiment,
                    condition, measurement, control:b))
```

	subject	sex	condition	measurement
## 1	S1	F	control	4.2
## 2	S2	M	control	5.9
## 3	S3	M	control	9.1
## 4	S5	F	control	2.1
## 5	S1	F	a	4.1
## 6	S2	M	a	7.2
## 7	S3	M	a	9.8
## 8	S5	F	a	23.5
## 9	S1	F	b	2.2

Long Data to Wide

- Use spread to move to a wide format

```
library(tidyr)
(data_wide = spread(data_long, condition, measurement))
```

```
##   subject sex    a    b control
## 1      S1   F  4.1  2.2      4.2
## 2      S2   M  7.2  6.8      5.9
## 3      S3   M  9.8 10.2      9.1
## 4      S5   F 23.5  5.2      2.1
```

On the Agenda

1 Data Management

- Motivation
- Sample Project
- Case Study: Housing Data
- Tidy Data

2 Regular Expressions

- Motivation
- Syntax
- Case Study: House Addresses
- Extracting and Formatting

3 Dates and Times

- Motivation
- System Information
- Operations on Time
- Date Formats
- Time Formats

4 Misc

- POSIXlt
- anytime
- lubridate

Regular Expressions

- **Regular Expression** or **regex** is a sequence of characters that defines a search pattern for a collection of strings.
- The idea sprouted from the notion of a **regular language** that was brought into existence by Kleene's theorem written by Stephen Cole Kleene.
- **Regex** is primarily used to:
 - ① search for patterns and,
 - ② replace patterns
- For it to function, programmers have adopted a set of grammatical statements to build patterns for strings.
- The grammar is available in just about every single programming language.

Regular Expression Usage Cases

- Validate Data Entry Fields
 - dates, e-mail address, credit card numbers
- Filter Text
 - key words or phrases in reviews, web server logs, reading config files
- Restructuring Text (Find and Replace)
 - mass change variable names, switching line endings
- Counting Occurrences
 - number of words, errors, or warnings

Relevant XKCD Comic



Figure 1: XKCD 208

Words of Wisdom for Regular Expressions

Some people, when confronted with a problem, think “I know, I’ll use regular expressions.” Now they have two problems.
— Jamie Zawinski in [alt.religion.emacs](#)

Note: Avoid using *regex* when parsers exist for tree structures (e.g. *html/dom*) to prevent edge cases from not being picked up!

Regular Expressions in R

Function	Description
<code>grep</code>	Returns a vector of the indices or values that match
<code>grep1</code>	Returns a logical vector indicating matches
<code>regexpr</code>	Returns the starting position of the first match
<code>gregexpr</code>	Returns the starting position of all matches
<code>sub</code>	Perform replacement of the first match
<code>gsub</code>	Perform replacement of the all matches

Regex Example - Finding IL (Concatenation)

- Consider the need to *filter* terms by whether or not they include **IL** (short for Illinois).

```
locs = c('Chicago, IL', 'San Francisco, CA', 'Springfield, IL',  
         'Detroit, MI', 'Urbana, IL', 'Tampa, FL')
```

```
grep('IL', locs)           # Obtain the Indices
```

```
## [1] 1 3 5
```

```
grep('IL', locs, value = TRUE) # Obtain the Names
```

```
## [1] "Chicago, IL"      "Springfield, IL" "Urbana, IL"
```

```
grepl('IL', locs)         # Obtain logical response
```

```
## [1] TRUE FALSE TRUE FALSE TRUE FALSE
```

Regex Example - Finding IL (Concatenation)

- Find the locations within the string
 - If missing location returns -1

```
regexpr('IL', locs) # Obtain the first instance in a word
```

```
## [1] 10 -1 14 -1 9 -1
## attr("match.length")
## [1] 2 -1 2 -1 2 -1
## attr("useBytes")
## [1] TRUE
```

```
gregexpr('IL', locs) # Obtain all instances in a word
```

```
## [[1]]
## [1] 10
## attr("match.length")
## [1] 2
## attr("useBytes")
## [1] TRUE
##
## [[2]]
```

Regex Example - Removing a (Concatenation)

- Remove instances of the letter a from the text

```
sub('a', "", locs) # Remove first instance
```

```
## [1] "Chicgo, IL"      "Sn Francisco, CA" "Springfield, IL"  
## [4] "Detroit, MI"     "Urbna, IL"       "Tmpa, FL"
```

```
gsub('a', "", locs) # Remove all instances
```

```
## [1] "Chicgo, IL"      "Sn Frncisco, CA" "Springfield, IL"  
## [4] "Detroit, MI"     "Urbn, IL"        "Tmp, FL"
```

Regex Operations

The operations for regex can be viewed in distinct groups.

- Foundational Operators
- Character Classes
- Quantifiers
- Anchors

Regex Lexicon - Round 1

Operation	Explanation	Symbol
Concatentation	Exact String	word
Wildcard	Any character	.
Union	Either character	
Closure	Match preceding character 0 or more	*
Parentheses	Matches a pattern group	()

Regex Lexicon - Examples for Round 1: Operators

Operation	Symbol	Example	Match	Failure
Concatentation	word	james	james	da yae
Wildcard	.	j.m.s	james/jomas	anmes
Union		jb am	jjb / jam	toad
Closure	*	cat*	catcat / ca	ma / at
Parentheses	()	(og)	dog / blog	dgs

Regex Lexicon - Round 2: Classes

- A **character class** is a *list* of characters enclosed between `[]` that matches *any single character* in that list.
- **Exception:** If the first character in the list is the caret `^` (e.g. `[^]`) than it matches any character not in the list.
- For example, the regular expression `[0123456789]` matches any *single* digit, and `[^abc]` matches anything except the characters a, b or c.

Operation	Explanation	Symbol
Class Matches	Match within specific classes	<code>[]</code>
Range	Match values within a range	<code>[-]</code>
Negations	Match any character not within	<code>[^]</code>

Regex Lexicon - Examples for Round 2: Classes

Operation	Symbol	Example	Match	Failure
Class Matches	[]	[abc]	toad / book	room / desk
Range	[-]	[a-zA-Z]	Funky / word	1234 / #\$\$
Negations	[^]	[^aeiou]	gst / wd	hi / orange

Regex Lexicon - Character Class Shortcuts

Operation	POSIX	Regex Default	Regex Shortcut
Space	<code>[:space:]</code>	<code>[]</code>	<code>\s</code>
Blank Chars	<code>[:blank:]</code>	<code>[\f\n\r\t\v]</code>	<code>\s</code>
Digits	<code>[:digit:]</code>	<code>[0-9]</code>	<code>\d</code>
Lower	<code>[:lower:]</code>	<code>[a-z]</code>	
Upper	<code>[:upper:]</code>	<code>[A-Z]</code>	
Punctuation	<code>[:punct:]</code>	Not listed	
Alphanumeric	<code>[:alnum:]</code>	<code>[a-zA-Z0-9_]</code>	<code>\w</code>
Alphabetic	<code>[:alpha:]</code>	<code>[a-zA-Z]</code>	

Notes:

- POSIX shortcuts adjust to locale vs. regex default.
- **All POSIX values must be wrapped in []**
 - e.g. `[:alpha:]` gives the alphabetic characters, for more see `?regex`

Regex Lexicon - Round 3: Quantifiers

Operation	Explanation	Symbol
At Most Once	Match preceding character 0 or 1 times	?
One or More	Match preceding character 1 or more	+
Exact Amount	Match exactly m occurrences	{ m }
At Least	Match at least m occurrences	{ m ,}
At Most	Match at most n occurrences	{0, n }
Between	Match $m \leq x \leq n$ occurrences	{ m , n }

Regex Lexicon - Examples for Round 3: Quantifiers

Operation	Symbol	Example	Match	Failure
At Most Once	?	a?	hat / car	no / yes
One or More	+	sh+	shoe / ship	hip / s
Exact Amount	{m}	[0-9]{3}	123-423 / 921	12-33 / 9
At Least	{m,}	[m]{2,}	yumm / mommy	mom / yum
At Most	{0,n}	(za){0,2}	pizza / zoids	Matches All
Between	{m,n}	[o]{1,2}	zoo / pod	dad / phil

Regex Lexicon - Round 4: Anchors (Other)

Operation	Explanation	Symbol
Beginning	Start at the beginning of the string	^
End	Start at the end of the string	\$

Regex Lexicon - Examples Round 4: Anchors (Other)

Operation	Symbol	Example	Match	Failure
Beginning	^	^hi	hiya / hillary	child / phi
End	\$	s\$	james / dogs	sam / sosa

Regex Escape sequences

Character	Sequence
Single Quote	\'
Double Quote	\"
Newline	\n
Carriage Return	\r
Tab	\t
New Page (Form Feed)	\f
Bell Character	\a

Regex Backreferences

- Regex has a nice feature that provides: **backreferences**
- Definition: **Backreferences** enable the retrieval of values within side the parantheses grouping option given by (). This is helpful when doing string replacement.
- Each () group is assigned a number n that can be referenced with `\n`.

```
(backref = c("xy", "xyz", "xayz", "x12z"))
```

```
## [1] "xy"    "xyz"   "xayz"  "x12z"
```

```
gsub("(xy)", "\\1stats", backref)
```

```
## [1] "xystats" "xystatsz" "xayz"      "x12z"
```

Interesting regex equivalences

Contained below are equivalent regex expressions

- e.g. *op1 == op2*

OP 1	OP 2	Equivalence Remark
<code>a+</code>	<code>aa*</code>	Replicates pattern
<code>a(b c)d</code>	<code>a[bc]d</code>	If single characters!
<code>a{0,1}</code>	<code>a?</code>	Optional pattern

Helpful Regex Tools

- [regex101.com](#): Great Regex tester
- [rubular](#): Also a great regex tester (specific to ruby)
- [regular-expressions](#): Lots of regex resources, but the color scheme is hard on the eyes.
- [txt2re](#): inline regex explanation

Extracting and Formatting regex Example

- Data may sometimes not come formatted to the necessary requirements.
- One variable might need to be split up among multiple variables.
- This is an ideal case for using **regular expressions** (?regex).

Extracting and Formatting

- In the housing data, the house address is stored as:

“[#### Street Name] , [City], [State] [Zipcode]”

- We would like to split this data into three new variables:
 - **St_Address**, Zipcode, and City.
- Note, we will want to use the lower case alphabet to prevent capitalization mismatches within alphabetical fields. (e.g. “harbor estates ln” \neq “Harbor Estates Ln”)

On the Agenda

1 Data Management

- Motivation
- Sample Project
- Case Study: Housing Data
- Tidy Data

2 Regular Expressions

- Motivation
- Syntax
- Case Study: House Addresses
- Extracting and Formatting

3 Dates and Times

- Motivation
- System Information
- Operations on Time
- Date Formats
- Time Formats

4 Misc

- POSIXlt
- anytime
- lubridate

Date and Time Formats

“The only reason for time is so that everything doesn’t happen at once.”

— *Albert Einstein*

- *R* has the ability to interface with time information.
- The interface, as we will see, may not be the best but it is highly versatile.
- This is important in a world that is going more and more global.

Date and Time Formats

```
Sys.Date()           # Returns a date as R's Date object
```

```
## [1] "2017-07-10"
```

```
Sys.time()           # Returns both date & time at current locale as POSIXct
```

```
## [1] "2017-07-10 14:43:32 CDT"
```

```
as.numeric(Sys.time()) # Seconds from UNIX Epoch (1970-01-01 00:00:00 UTC)
```

```
## [1] 1499715812
```

```
Sys.timezone()       # Time zone at current location
```

```
## [1] "America/Chicago"
```

Date and Time Formats - Failure of characters

- Frequently, dates and times will be given as characters within a `data.frame`
- Having dates as characters impedes ones ability to be able to use the time information in an analysis
 - For example: How long did it take for the help desk call to be completed?

Bad Time Differencing:

```
time1 = "2016-06-28 10:25:44 CDT" # UNIX Time Stamp
time2 = "2016-06-28 15:25:44 CDT" # UNIX Time Stamp
time2 - time1
# Error in time2 - time1 : non-numeric argument to
# binary operator
```

Date and Time Formats - Time Operations

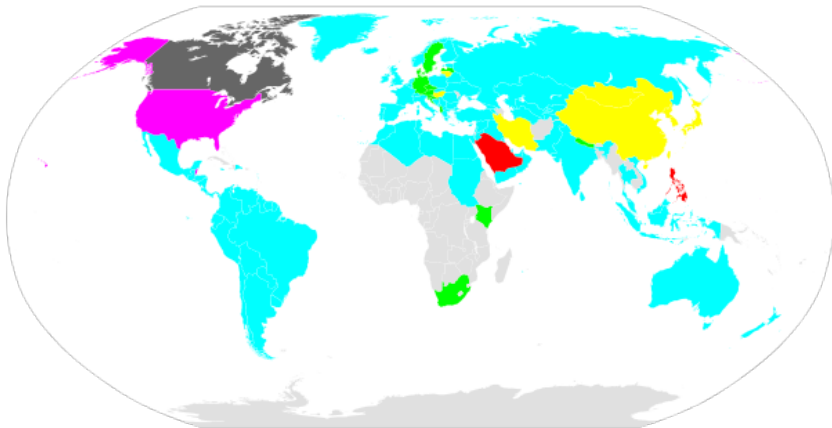
- Performing time operations requires that both dates are given as POSIXct object in R.

```
time1 = as.POSIXct("2016-06-28 10:25:44")  
time2 = as.POSIXct("2016-06-28 15:25:44") # +5 Hours  
time2 - time1
```

Time difference of 5 hours

Note: Default format for POSIXct is %Y-%m-%d %H:%M:%S

The Date Format Around the World



Color	Date Format	Main Region	Population (Millions)
Cyan	DD/MM/YYYY	Australia, Russia	3295
Magenta	MM/DD/YYYY	USA, Canada	360
Yellow	DD/MM/YYYY	China, India	1400
Red	DD/MM/YYYY	Middle East	300
Green	DD/MM/YYYY	South Africa, parts of Europe	100
Grey	Unknown / Multiple	Africa, South America, Oceania	1000

Formats for Working with Dates

Format	Description	Example
%a	Abbreviated weekday name in the current locale	Mon
%A	Full weekday name in the current locale	Monday
%b	Abbreviated month name in the current locale	Jul
%B	Full month name in the current locale	July
%m	Month number (01-12)	07
%d	Day of the month as decimal number (01-31)	10
%e	Day of the month as decimal number (1-31)	10
%y	Year without century (00-99)	17
%Y	Year including century	2017

For more, see `?strptime`

Formating Non-Standard Dates

```
(yyyy_mm_dd = as.POSIXct("2016-06-28",  
                           format = "%Y-%m-%e"))
```

```
## [1] "2016-06-28 CDT"
```

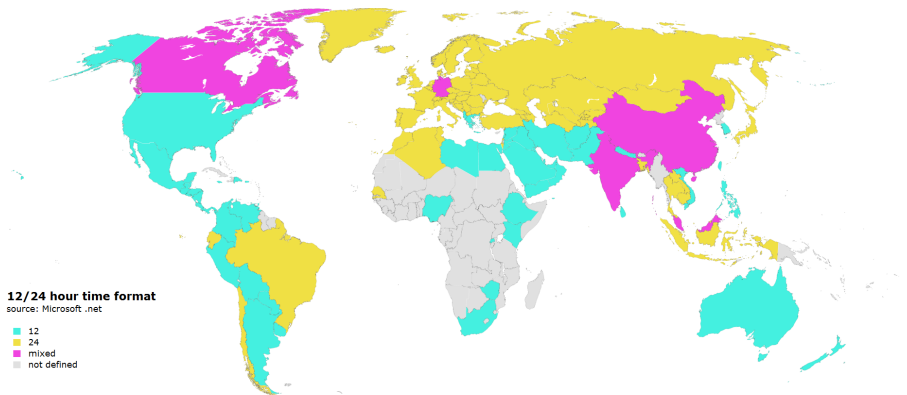
```
(dd_mm_yy = as.POSIXct("28/06/16",  
                        format = "%e/%m/%y"))
```

```
## [1] "2016-06-28 CDT"
```

```
(mon_dd_yyyy = as.POSIXct("Jun 28, 2016",  
                           format = "%b %e, %Y"))
```

```
## [1] "2016-06-28 CDT"
```

Time Format Used Around the World



Formats for Working with Times

Format	Description	Example
%S	Second as integer (00–61)	32
%OS	Second as decimal number (00-60.99)	32
%M	Minute as decimal number (00–59)	43
%H	Hours as decimal number (00–23)	14
%I	Hours as decimal number (01–12)	02
%p	AM/PM indicator in the locale	PM
%z	Signed offset in hours and minutes from UTC	-0500
%Z	Time zone abbreviation as a character string	CDT

For more, see `?strptime`

Formating Non-Standard Times

```
(h_m = as.POSIXct("11:38",  
                  format = "%H:%M"))
```

```
## [1] "2017-07-10 11:38:00 CDT"
```

```
(h_am = as.POSIXct("11 AM",  
                   format = "%I %p"))
```

```
## [1] "2017-07-10 11:00:00 CDT"
```

```
(h_m_s_z = as.POSIXct("11:38:22", # Chop off the TZ  
                      format = "%H:%M:%S",  
                      tz = "America/New_York"))
```

```
## [1] "2017-07-10 11:38:22 EDT"
```

Time Zone Notes

- *R* makes use of time zones via `tz` parameter.
- The accepted values of `tz` depend on the location.
 - CST is given with `"CST6CDT"` or `"America/Chicago"`
- For supported locations and time zones use:
 - In *R*: `OlsonNames()`
 - Alternatively, try in *R*: `system("cat $R_HOME/share/zoneinfo/zone.tab")`
- These locations are given by [Internet Assigned Numbers Authority \(IANA\)](#)
 - [List of tz database time zones \(Wikipedia\)](#)
 - [IANA TZ Data \(2016e\)](#)

Specifics on POSIXct

- POSIXct: Stores time as seconds since UNIX epoch on 1970-01-01 00:00:00
 - Unix & tidyverse preferred format.

```
# POSIXct output
(origin = as.POSIXct("1970-01-01 00:00:00",
                      format = "%Y-%m-%d %H:%M:%S",
                      tz = "UTC"))
```

```
## [1] "1970-01-01 UTC"
```

```
as.numeric(origin)      # At epoch
```

```
## [1] 0
```

```
as.numeric(Sys.time()) # Right now
```

```
## [1] 1499715813
```

On the Agenda

1 Data Management

- Motivation
- Sample Project
- Case Study: Housing Data
- Tidy Data

2 Regular Expressions

- Motivation
- Syntax
- Case Study: House Addresses
- Extracting and Formatting

3 Dates and Times

- Motivation
- System Information
- Operations on Time
- Date Formats
- Time Formats

4 Misc

- `POSIXlt`
- `anytime`
- `lubridate`

The “other” time object: POSIXlt

- POSIXlt: Stores a list of day, month, year, hour, minute, second, and so on.
 - It is **slower** than POSIXct and has **zero support** in the tidyverse.
 - **Warning:** POSIXlt will be returned if you use `strptime()`
 - Always convert POSIXlt to POSIXct using `as.POSIXct()!!!`

```
# POSIXlt output
```

```
posixlt = as.POSIXlt(Sys.time(),  
                      format = "%Y-%m-%d %H:%M:%S",  
                      tz = "America/Chicago")
```

```
# Convert to POSIXct
```

```
posixct = as.POSIXct(posixlt)  
posixct
```

```
## [1] "2017-07-10 14:43:33 CDT"
```

POSIXlt - List Values

```
posixlt$sec    # Seconds 0-61
```

```
## [1] 33.49249
```

```
posixlt$min    # Minutes 0-59
```

```
## [1] 43
```

```
posixlt$hour   # Hour 0-23
```

```
## [1] 14
```

```
posixlt$mday   # Day of the Month 1-31
```

```
## [1] 10
```

```
posixlt$mon    # Months after the first of the year 0-11
```

```
## [1] 6
```

```
posixlt$year   # Years since 1900.
```

```
## [1] 117
```

anytime

- `anytime` by Dirk Eddelbuettel seeks to solve the need of remembering date and time formats.
- Main advantage: Only one function `anytime()` that autodetects the appropriate format and imports it correctly as a `POSIXct` object.

```
library(anytime)
Sys.setenv(TZ=anytime:::getTZ()) ## helper function to try to get TZ

anytime(c("2017-Jul-10 10:11:12", "Jul/10/2017 10:11:12", "Jul-10-2017 10:11:12"))
```

```
## [1] "2017-07-10 10:11:12 CDT" "2017-07-10 10:11:12 CDT"
## [3] "2017-07-10 10:11:12 CDT"
```

```
anytime(c("Mon Jul 10 10:11:12 2016", "Mon Jul 10 10:11:12.345678 2017"))
```

```
## [1] "2016-07-10 10:11:12 CDT" "2017-07-10 10:11:12 CDT"
```

lubridate - Dates Made Easy

- lubridate by Garret Grolemund, Hadley Wickham, and Gang contains **many** helper functions to write the correct parse syntax e.g.

```
library(lubridate)
ymd("20160628")
```

```
## [1] "2016-06-28"
```

```
interval(mdy("06-28-2016"), dmy("29/06/2016"))
```

```
## [1] 2016-06-28 UTC--2016-06-29 UTC
```

For more, please read the [Lubridate vignette](#) on the [CRAN](#)

Summary

- To analyze time, you must have it in a POSIXct object
 - Avoid POSIXlt like the plague.
- Date and Timestamps differ greatly around the world.
- Lots of options outside of base R functions exist.