



Journal of Statistical Software

MMMMMM YYYY, Volume VV, Issue II.

<http://www.jstatsoft.org/>

Dates and Times Made Easy with lubridate

Garrett Grolemund
Rice University

Hadley Wickham
Rice University

Abstract

This paper presents the lubridate package for R (R Development Core Team 2010), which facilitates working with dates and times. Date-times create various technical problems for the data analyst. The paper highlights these problems and offers practical advice on how to solve them using **lubridate**. The paper also introduces a conceptual framework for arithmetic with date-times.

Keywords: dates, times, time zones, daylight savings time, R.

1. Introduction

Date-time data can be frustrating to work with. Dates come in many different formats, which makes recognizing and parsing them a challenge. Will our program handle the format that we have? If it does, we still face problems specific to date-times. How can we easily extract components of the date-times, such as years, months, or seconds? How can we switch between time zones, or compare times from places that use daylight savings time (DST) with times from places that do not? Date-times create even more complications when we try to do arithmetic with them. Conventions such as leap years and DST make it unclear what we mean by “one day from now” or “exactly two years away.” Even leap seconds can disrupt a seemingly simple calculation. This complexity affects other tasks too, such as constructing sensible tick marks for plotting date-time data.

While Base R (R Development Core Team 2010) handles some of these problems, the syntax it uses can be confusing and difficult to remember. Moreover, the correct R code often changes depending on the class of date-time object being used. **lubridate** acknowledges these problems and makes it easier to work with date-time data in R. **lubridate** will enhance a user’s experience for any analysis that includes date-time data. It also provides tools for manipulating date-times in novel but useful ways. Specifically, **lubridate** helps users:

- Identify and parse date-time data, Section 3.
- Extract and modify components of a date-time, such as years, months, days, hours, minutes, and seconds, Section 4.
- Perform accurate math on date-times, Section 5.
- Handle time zones and Daylight Savings Time, Sections 7, and 8.

lubridate uses an intuitive user interface inspired by the date libraries of object oriented programming languages. **lubridate** methods are compatible with a wide-range of common date-time objects. These include character strings, POSIXct, POSIXlt, Date, **chron** (James and Hornik 2010), **fCalendar** (Wuertz, with contributions from Martin Maechler, Byers *et al.* 2010), **zoo** (Zeileis and Grothendieck 2005), **xts** (Ryan and Ulrich 2010), **its** (Portfolio and Risk Advisory Group, Commerzbank Securities 2009), **tis** (Hallman 2010), **timeSeries** (Wuertz and Chalabi 2010), **fts** (Armstrong 2009), and **tseries** (Trapletti and Hornik 2009) objects.

Note that **lubridate** overrides the + and - methods for POSIXt, Date, and difftime objects in base R. This allows users to perform simple arithmetic on date-time objects with the new time span classes introduced by **lubridate**, but it does not alter the way R implements addition and subtraction for non-**lubridate** objects.

lubridate introduces four new object classes based on the Java based JODA time project (Colebourne and O'Neill 2006). JODA time introduces a conceptual model of the different ways to measure time spans. Section 5 describes this model and explains how **lubridate** uses it to perform easy and accurate arithmetic with dates.

This paper demonstrates the convenient tools provided in the **lubridate** package and ends with a case study, which uses **lubridate** in a real life example. This paper describes **lubridate** 0.1, which can be downloaded from CRAN. Development versions can be found at <http://github.com/hadley/lubridate>.

2. Motivation

To see how **lubridate** simplifies things, consider a common scenario. Given a character string, we would like to read it in as a date-time, extract the month, and change it to February (i.e. 2). On the left are the base R methods we'd use for these three tasks. On the right are the **lubridate** methods.

<pre>date <- as.POSIXct("01-01-2010", format = "%d-%m-%Y", tz = "UTC") as.numeric(format(date, "%m")) # or as.POSIXlt(date)\$month + 1 date <- as.POSIXct(format(date, "%Y-2-%d"), tz = "UTC")</pre>	<pre>date <- dmy("01-01-2010") month(date) month(date) <- 2</pre>
--	---

Now let's go a step further. We'll go back in time one day and display our new date in the Greenwich Meridian time zone (GMT). Again, base R methods are shown on the left, **lubridate** methods on the right.

<pre>date <- seq(date, length = 2, by = "-1 day")[2] as.POSIXct(format(as.POSIXct(date), tz = "UTC"), tz = "GMT")</pre>	<pre>date <- date - days(1) with_tz(date, "GMT")</pre>
---	--

lubridate makes basic date-time manipulations much more straightforward. Plus, the same **lubridate** methods work for most of the popular date-time object classes (`Date`, `POSIXt`, `chron`, etc.).

Table 1 provides a more complete comparison between **lubridate** methods and base R methods. It shows how **lubridate** can simplify each of the common date time tasks presented in R News Volume 4/1, June 2004 (Grothendieck and Petzoldt 2004). It also provides a useful summary of **lubridate** methods.

3. Parsing date-times

We can read dates into R using the `ymd()` series of functions provided by **lubridate**. The letters y, m, and d correspond to the year, month, and day elements of a date-time. To read in a date, choose the function name that matches the order of elements in your date-time object. For example, in the following date the month element comes first, followed by the day and then the year. So we'd use the `mdy()` function:

```
R> mdy("12-01-2010")
[1] "2010-12-01 UTC"
```

The same character string can be parsed as January 12, 2001 by reversing the month and day element with `dmy()`.

```
R> dmy("12-01-2010")
[1] "2010-01-12 UTC"
```

The `ydm()` series of functions can also parse vectors of dates.

```
R> dmy(c("31.12.2010", "01.01.2011"))
[1] "2010-12-31 UTC" "2011-01-01 UTC"
```

These functions create a `POSIXct` date-time object that matches the date described by the character string. The functions automatically recognize the separators commonly used to record dates. These include: “-”, “/”, “.”, and “” (i.e., no separator). When a `ymd()` function is applied to a vector of dates, **lubridate** will assume that all of the dates have the same order and the same separators. `ymd()` type functions also exist for times recorded with hours,

Task	Lubridate	Date	POSIXct
now (system time zone)	now()	Sys.Date()	Sys.time()
now (GMT)	now("GMT")	structure(0, class = "Date")	structure(0, class = c("POSIXt", "POSIXct"))
origin	origin	structure(floor(x), class = "Date")	structure(x*24*60*60, class=c("POSIXt", "POSIXct"))
x days since origin	origin + days(x)	date + 1	seq(date, length = 2, by = "day") [2]
next day	date + days(1)	date + 1	seq(date, length = 2, by = "-1 day") [2]
previous day	date - days(1)	date - 1	seq(date, length = 2, by = "-1 day") [2]
DST and Time Zones			
x days since date	date + ddays(x)	date + floor(x)	seq(date, length = 2, by = paste(x, "day")) [2]
(day exactly 24 hours)	date + days(x)		seq(date, length = 2, by = paste(x, "DSTday")) [2]
(allowing for DST)	with_tz(date, "TZ")		as.POSIXct(format(as.POSIXct(date), tz = "TZ"), tz = "TZ")
display date in new time zone			
keep clock time, replace time zone	force_tz(date, tz = "TZ")		
Exploring			
sequence	date + c(0:9) * days(1)	seq(date, length = 10, by = "day")	seq(date, length = 10, by = "DSTday")
every 2nd week	date + c(0:2) * weeks(2)	seq(date, length = 3, by = "2 week")	seq(date, length = 3, by = "week")
first day of month	floor_date(date, "month")	as.Date(format(date, "%Y-%m-01"))	as.POSIXct(format(date, "%Y-%m-01"))
round to nearest first of month	round_date(date, "month")		
extract year value	year(date)	as.numeric(format(date, "%Y"))	as.numeric(format(date, "%Y"))
change year value	year(date) <- z	as.Date(format(date, "%z-%m-%d"))	as.POSIXct(format(date, "%z-%m-%d"))
day of week	yday(date) # Sun = 1	as.numeric(format(date, "%w")) # Sun = 0	as.numeric(format(date, "%w")) # Sun = 0
day of year	yday(date)	as.numeric(format(date, "%j"))	as.numeric(format(date, "%j"))
express as decimal of year	decimal_date(date)		
Parse			
z = "1970-10-15"	ymd(z)	as.Date(z)	as.POSIXct(z)
z = "10/15/1970"	mdy(z)	as.Date(z, "%m/%d/%Y")	as.POSIXct(strptime(z, "%m/%d/%Y"))
z = 15101970	dmy(z)	as.Date(as.character(z), format = "%d%ym%Y")	as.POSIXct(as.character(z), tz = "GMT", format = "%d%ym%Y")
Durations Comparison			
Duration			
1 second	seconds(1)	base::R	as.difftime(1, unit = "secs")
5 days, 3 hours and - 1 minute	new_duration(day = 5, hour = 3, minute = -1)		as.difftime(60 * 24 * 5 + 60 * 3 - 1, unit = "mins")
1 month	months(1)		# Time difference of 7379 mins
1 year	years(1)		

Table 1: Comparison Table (adapted from R News Volume 4/1, June 2004).

Order of elements in date-time	Parse function
year, month, day	<code>ymd()</code>
year, day, month	<code>ydm()</code>
month, day, year	<code>mdy()</code>
day, month, year	<code>dmy()</code>
hour, minute	<code>hm()</code>
hour, minute, second	<code>hms()</code>
year, month, day, hour, minute, second	<code>ymd_hms()</code>

Table 2: Parse function names are based on the order that years, months, and days appear within the dates to be parsed.

minutes, and seconds. These functions make it simple to parse any date-time object that can be converted to a character string. See Table 2 for a complete list of `ymd()` type parsing functions.

4. Manipulating date-times

Every date-time is a combination of different elements, each with its own value. For example, most date-times include a year value, a month value, a day value and so on. Together these elements specify the exact moment that the date-time refers to. We can easily extract each element of a date-time with the accessor function that has its name, as shown in Table 3. For example, if we save the current system time

```
R> date <- now()
[1] "2010-02-25 09:51:48 CST"
```

we can extract each of its elements. Note that this was the system time when this example was written. `now()` will return a different date-time each time it is used.

```
R> year(date)
[1] 2010
```

```
R> minute(date)
[1] 51
```

For the month and weekday elements (`wday`), we can also specify whether we want to extract the numerical value of the element, an abbreviation of the name of the month or weekday, or the full name. For example,

```
R> month(date)
[1] 2
```

```
R> month(date, label = TRUE)
```

Date component	Accessor
Year	<code>year()</code>
Month	<code>month()</code>
Week	<code>week()</code>
Day of year	<code>yday()</code>
Day of month	<code>mday()</code>
Day of week	<code>wday()</code>
Hour	<code>hour()</code>
Minute	<code>minute()</code>
Second	<code>second()</code>
Time Zone	<code>tz()</code>

Table 3: Each date-time element can be extracted with its own accessor function.

```
[1] Feb
```

```
R> month(date, label = TRUE, abbr = FALSE)
```

```
[1] February
```

```
R> wday(date, label = TRUE, abbr = FALSE)
```

```
[1] Thursday
```

We can also use any of the accessor functions to set the value of an element. This would also change the moment that the date-time refers to. For example,

```
R> day(date) <- 5
```

```
[1] "2010-02-05 09:51:48 CST"
```

changes our date to the fifth day of the month. We can also set the elements to more complicated values, e.g.

```
R> dates <- ymd_hms("2010-01-01 01:00:00", "2010-01-01 01:30:00")
```

```
R> minute(dates) <- mean(minute(dates))
```

```
[1] "2010-01-01 01:15:00 UTC" "2010-01-01 01:15:00 UTC"
```

Note that if we set an element to a larger value than it supports, the difference will roll over into the next higher element. For example,

```
R> day(date) <- 30
```

```
[1] "2010-03-02 09:51:48 CST"
```

We can use this to find the last day of a month:

```
R> day(date) <- 1
R> month(date) <- month(date) + 1
R> day(date) <- day(date) - 1
[1] "2010-03-31 09:51:48 CDT"
```

Lubridate also provides an update method for date-times. This is useful if you want to change multiple attributes at once or would like to create a modified copy, instead of transforming in place.

```
R> update(date, year = 2010, month = 1, day = 1)
[1] "2010-01-01 09:51:48 CST"
```

Finally, we can also change dates by adding or subtracting units of time from them. For example, the methods below produce the same result.

```
R> hour(date) <- 12
[1] "2010-02-25 12:51:48 CST"
R> date + hours(3)
[1] "2010-02-25 12:51:48 CST"
```

Notice that `hours()` (plural) is not the same function as `hour()` (singular). `hours()` creates a new object that can be added or subtracted to a date-time. These objects are discussed in the next section.

5. Math with date-times

Arithmetic with date-times is more complicated than arithmetic with numbers, but it can be done accurately and easily with **lubridate**. What complicates math with date-times? Clock times are periodically re-calibrated to reflect astronomical conditions, such as the hour of daylight, or the Earth's tilt on its axis relative to the sun. We know these re-calibrations as daylight savings time, leap years, and leap seconds. Consider how one of these conventions might complicate a simple math operation. If today were January 1st, 2010 and we wished to know what day it would be one year from now, we could simply add 1 to the years element of our date.

January 1st, 2010 + 1 year = January 1st, 2011

Alternatively, we could add 365 to the days element of our date because a year is equivalent to 365 days.

January 1st, 2010 + 365 days = January 1st, 2011

Troubles arise if we try the same for January 1st, 2012. 2012 is a leap year, which means it has an extra day. Our two approaches above now give us different answers because the length

of a year has changed.

January 1st, 2012 + 1 year = January 1st, 2013

January 1st, 2012 + 365 days = December 31st, 2012

At different moments in time, the lengths of months, weeks, days, hours, and even minutes will also vary. We can consider these to be *relative* units of time; their length is relative to when they occur. In contrast, seconds always have a consistent length. Hence, seconds are *exact* units of time.

Researchers may be interested in exact lengths, relative lengths, or both. For example, the speed of a physical object is most precisely measured in exact lengths. The opening bell of the stock market is more easily modeled with relative lengths.

lubridate allows math with both relative and exact units by introducing four new time related objects. These are *instants*, *intervals*, *durations*, and *periods*, terminology borrowed from the JODA time project (Colebourne and O'Neill 2006).

5.1. Instants

An instants is a specific moment in time, such as January 1st, 2012. We create an instant each time we parse a date into R.

```
R> start_2012 <- ymd_hms("2012-01-01 12:00:00")
```

lubridate does not create instant class objects. Instead, it recognizes any date-time object that refers to a moment of time as an instant. We can test if an object is an instant by using `is.instant()`. For example,

```
R> is.instant(364)
[1] FALSE
```

```
R> is.instant(start_2012)
[1] TRUE
```

We can easily round instants to the nearest minute, hour, month, etc. using `floor_date()`, `ceiling_date()`, and `round_date()`. For example,

```
R> round_date(start_2012, "day")
[1] "2012-01-02 UTC"
```

We can also capture the current time as an instant with `now()`, and the current day with `today()`.

5.2. Intervals

Intervals, durations, and periods are all ways of recording time spans. Of these, intervals are the most simple. An interval is a period of time that occurs between two specific instants. The

length of an interval is never ambiguous, because we know when it occurs. We can calculate the exact lengths of any unit of time that occurs during it.

We can create interval objects by subtracting two instants or by using the command `new_interval()`.

```
R> start_2011 <- ymd_hms("2011-01-01 12:00:00")
R> start_2010 <- ymd_hms("2010-01-01 12:00:00")
R> span <- start_2011 - start_2010
[1] 365 days beginning at 2010-01-01 12:00:00
```

Unfortunately, since intervals are anchored to their start and end dates, they are not very useful for date-time math. It only makes sense to add an interval to its start date or to subtract it from its end date.

```
R> start_2010 + span
[1] "2011-01-01 12:00:00 UTC"
```

5.3. Durations

If we remove the start and end dates from an interval, we will have a generic time span that we can add to any date. But how should we measure this length of time? If we record the time span in seconds, it will have an exact length since seconds always have the same length. We call such time spans *durations*. Alternatively, we can record the time span in larger units, such as minutes or years. Since the length of these units varies over time, the exact length of the time span will depend on when it begins. These non-exact time spans are called *periods* and will be discussed in the next section.

The length of a duration is invariant to leap years, leap seconds, and daylight savings time because durations are measured in seconds. Hence, durations have consistent lengths and can be easily compared to other durations. Durations are the appropriate object to use when comparing time based attributes, such as speeds, rates, and lifetimes.

lubridate uses the `difftime` class from base R for durations. Additional `difftime` methods have been created to facilitate this.

For large durations, it becomes inconvenient to describe the length in seconds. For example, not many people would recognize 31536000 seconds as the length of a standard year. For this reason, **lubridate** uses standard time units to display durations. However, these units are only estimates given for convenience. The underlying object is always recorded as a fixed number of seconds. Estimated units are created using the following relationships. A minute is 60 seconds, an hour 3600 seconds, a day 86400, a week 604800, and a year 31536000. Month units are not used because they are so variable.

Duration objects can be easily created with the helper functions `dyears()`, `dweeks()`, `ddays()`, `dhours()`, `dminutes()`, and `dseconds()`. The `d` in the title stands for duration and distinguishes these objects from period objects, which are discussed in Section 5.4. Each object creates a duration in seconds using the estimated relationships given above. The argument of each function is the number of estimated units we wish to include in the duration. For example,

```
R> eminutes(1)
Time difference of 1 mins
```

```
R> esconds(60)
Time difference of 1 mins # 60 seconds = 1 estimated minute
```

```
R> eminutes(2)
Time difference of 2 mins
```

```
R> c(1:3) * ehours(1)
Time difference in hours
[1] 1 2 3
```

Durations can be added and subtracted to any instant object. For example,

```
R> start_2011 + eyears(1)
[1] "2012-01-01 12:00:00 UTC"
```

```
R> start_2012 <- ymd_hms("2012-01-01 12:00:00")
R> start_2012 + eyears(1)
[1] "2012-12-31 12:00:00 UTC"
```

Durations can also be added to or subtracted from intervals and other durations. For example,

```
R> eweeks(1) + edays(6) + ehours(2) + eminutes(1.5) + esconds(3)
Time difference of 1.869201 weeks
```

We can also create durations from intervals using `as.duration()`.

```
R> as.duration(span)
Time difference of 52.14286 weeks # 1 year
```

5.4. Periods

Periods record a time span in units larger than seconds, such as years, months, weeks, days, hours, and minutes. For convenience, we can also create a period that only uses seconds, but such a period would have the same properties as a duration. We construct periods with the helper functions `years()`, `months()`, `weeks()`, `days()`, `hours()`, `minutes()`, and `seconds()`.

```
R> months(3)
[1] 3 months
```

```
R> months(3) + days(2)
```

	Instant	Interval	Duration	Period
Instant	N/A	instant	instant	instant
Interval	instant	N/A	interval	interval
Duration	instant	interval	duration	N/A
Period	instant	interval	N/A	period

Table 4: Adding two date-time objects will create the above type of object.

```
[1] 3 months and 2 days
```

These functions do not contain an `e` in their name, because they are no longer estimates. For example, `months(2)` always has the length of two months even though the length of two months will change depending on when the period begins. For this reason, we can not compute exactly how long a period will be in seconds until we know when it occurs. However, we can still do date-time math with periods. When we add or subtract a period to an instant, the period becomes anchored to the instant. The instant tells us when the period occurs, which allows us to calculate its exact length in seconds.

In other words, we can use periods to accurately model clock times without knowing when events such as leap seconds, leap days, and DST changes occur.

```
R> start_2012 + years(1)
[1] "2013-01-01 12:00:00 UTC"
```

vs.

```
R> start_2012 + eyears(1)
[1] "2012-12-31 12:00:00 UTC"
```

We can also convert intervals to period objects with the function `as.period()`.

```
R> as.period(span)
[1] 1 year
```

Periods can be added to instants, intervals, and other periods, but not to durations.

In summary, math with date-times involves four types of objects: instants, intervals, durations, and periods. Table 4 describes which objects can be added to each other and what type of object will result.

6. Rounding dates

Like numbers, date-times occur in order. This allows date-times to be rounded and `lubridate` provides three methods that help perform this rounding: `round_date()`, `floor_date()`, and `ceiling_date()`. The first argument of each function is the date-time or date-times to be rounded. The second argument is the unit to round to. For example, we could round April 20, 2010 to the nearest day

```
R> april20 <- ymd_hms("2010-04-20 11:33:29")
R> round_date(april20, "day")
[1] "2010-04-20 UTC"
```

or the nearest month.

```
R> round_date(april20, "month")
[1] "2010-05-01 UTC"
```

Notice that rounding a date-time to a unit sets the date to the start of that unit (e.g, `round_date(april20, "day")` sets the hours, minutes, and seconds information to 00).

`ceiling_date()` provides another simple way to find the last day of a month. Ceiling the date to the next month and then subtract a day.

```
R> ceiling_date(april20, "month") - days(1)
[1] "2010-04-30 UTC"
```

7. Time zones

Time zones give multiple names to the same instant. For example, “2010-03-26 11:53:24 CDT” and “2010-03-26 12:53:24 EDT” both describe the same instant. The first shows how the instant is labeled in the United States’ central time zone (CDT). The second shows how the same instant is labelled in the United States’ eastern time zone (EDT). Time zones complicate date-time data but are useful for mapping clock time to local daylight conditions. When working with instants, it is standard to give the clock time as it appears in the Coordinated Universal time zone (UTC). This saves calculations but can be annoying if your computer insists on translating times to your current time zone. It may also be inconvenient to discuss clock times that occur in a place unrelated to the data.

lubridate eases the frustration caused by time zones in two ways. We can change the time zone in which an instant is displayed by using the function `with_tz()`. This changes how the clock time is displayed, but not the instant that is referred to. For example,

```
R> date
[1] "2010-01-01 09:51:48 CST"
R> with_tz(date, "UTC")
[1] "2010-01-01 15:51:48 UTC"
```

`force_tz()` does the opposite of `with_tz()`: it changes the instant that is displayed, but the clock time remains the same. For example, the code below moves us to a new instant that occurs 6 hours earlier.

```
R> date
```

```
[1] "2010-01-01 09:51:48 CST"
R> force_tz(date, "UTC")
[1] "2010-01-01 09:51:48 UTC"
```

8. Daylight savings time

In many parts of the world, the official clock time springs forward by one hour in the spring and falls back one hour in the fall. For example, in Chicago, Illinois a change in daylight savings time occurred at 2:00AM on March 14, 2010. The last instant to occur before this change was 2010-03-14 01:59:59 CST.

```
R> dst_time <- ymd_hms("2010-03-14 01:59:59")
R> dst_time <- force_tz(dst_time, "America/Chicago")
[1] "2010-03-14 01:59:59 CST"
```

One second later, Chicago clock times read

```
R> dst_time + esconds(1)
[1] "2010-03-14 03:00:00 CDT"
```

It seems that we gained an extra hour during that second, which is how daylight savings time works. We can completely avoid the changes in clock time caused by daylight savings times by working with periods instead of durations. For example,

```
R> dst_time + hours(2)
[1] "2010-03-14 03:59:59 CDT"
```

displays the clock time that usually occurs two hours after 1:59:59 AM. When using periods, we do not have to track DST changes because they will not affect our calculations. Adding a duration would give us the actual clock time that appeared exactly two hours later on March 14, 2010.

```
R> dst_time + ehours(2)
[1] "2010-03-14 04:59:59 CDT"
```

If we ever inadvertently try to create an instant that occurs in the one hour “gap” between 2010-03-14 01:59:59 CST and 2010-03-14 03:00:00 CDT, **lubridate** will return NA since such times are not available.

We can also avoid the complications created by daylight savings time by keeping our date-times in a time zone such as “UTC”, which does not adopt daylight savings hours.

9. Case study 1

The next two sections will work through some techniques using **lubridate**. First, we will use **lubridate** to calculate the dates of holidays. Then we’ll use **lubridate** to explore an example

data set (`lakers`).

9.1. Thanksgiving

Some holidays, such as Thanksgiving (U.S.) and Memorial Day (U.S.) do not occur on fixed dates. Instead, they are celebrated according to a common rule. For example, Thanksgiving is celebrated on the fourth Thursday of November. To calculate when Thanksgiving will be held in 2010, we can start with the first day of 2010.

```
R> date <- ymd("2010-01-01")
[1] "2010-01-01 UTC"
```

We can then add 10 months to our date, or directly set the date to November.

```
R> month(date) <- 11
[1] "2010-11-01 UTC"
```

We check which day of the week November 1st is.

```
R> wday(date, label = T, abbr = F)
[1] Monday
```

This implies November 4th will be the first Thursday of November.

```
R> date <- date + days(3)
[1] "2010-11-04 UTC"
R> wday(date, label = T, abbr = F)
[1] Thursday
```

Next, we add three weeks to get to the fourth Thursday in November.

```
R> date + weeks(3)
[1] "2010-11-25 UTC"
```

9.2. Memorial Day

Memorial day also occurs according to a rule; it always occurs on the last Monday of May. To calculate the date of Memorial day, we can again start with the first of the year.

```
R> date <- ymd("2010-01-01")
[1] "2010-01-01 UTC"
```

Next, we set the month to May.

```
R> month(date) <- 5
[1] "2010-05-01 UTC"
```

Now our holiday occurs in relation to the last day of the month instead of the first. We find the last day of the month by rounding up to the next month using `ceiling_date()` and then subtracting a day.

```
R> date <- ceiling_date(date, "month") - days(1)
[1] "2010-05-31 UTC"
```

We can then check which day of the week May 31st falls on. It happens to be a Monday, so we are done. If May 31st had been another day of the week, we could've subtracted an appropriate number of days to get the last Monday of May.

```
R> wday(date, label = T, abbr = F)
[1] Monday
```

10. Case study 2

The `lakers` data set contains play by play statistics of every major league basketball game played by the Los Angeles Lakers during the 2008-2009 season. This data is from <http://www.basketballgeek.com/downloads/> (Parker 2010) and comes with the `lubridate` package. We will explore the distribution of Lakers' games throughout the year as well as the distribution of plays within Lakers' games. We choose to use the `ggplot2` (Wickham 2009) package to create our graphs.

The `lakers` data set comes with a `date` variable which records the date of each game. Using the `str()` command, we see that R recognizes the dates as integers.

```
str(lakers$date)
int [1:34624] 20081028 20081028 20081028 ...
```

Before we can work with them as dates, we must parse them into R as date-time objects. The dates appear to be arranged with their year element first, followed by the month element, and then the day element. Hence, we should use the `ymd()` parsing function.

```
R> lakers$date <- ymd(lakers$date)
str(lakers$date)
POSIXct[1:34624], format: "2008-10-28" "2008-10-28" ...
```

R now recognizes the dates as POSIXct date-time objects. It will now treat them as date-times in any functions that have POSIXct specific methods. For example, if we plot the occurrences of home and away games throughout the season, our x axis will display date-time

information for the tick marks (Figure 1).

```
R> qplot(date, 0, data = lakers, colour = game_type)
```

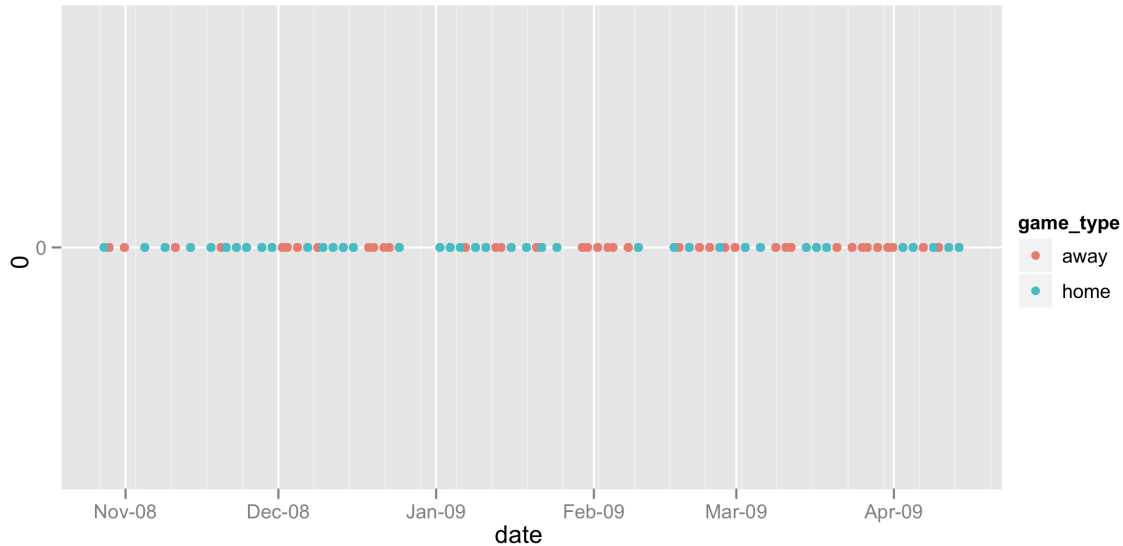


Figure 1: Occurrences of home and away games for the 2008-2009 season.

Figure 1 shows that games are played continuously throughout the season with a few short breaks. The frequency of games seems lower at the start of the season and games appear to be grouped into clusters of home games and away games. The tick marks and breaks on the x axis are automatically generated by the **lubridate** method `pretty.dates()`.

Next we'll examine how Lakers games are distributed throughout the week. We use the `wday()` command to extract the day of the week of each date.

```
R> qplot(wday(date, label = T), data = lakers, geom = "histogram")
```

The frequency of basketball games varies throughout the week (Figure 2). Surprisingly, the highest number of games are played on Tuesdays.

Now let's look at the games themselves. In particular, let's look at the distribution of plays throughout the game. The `lakers` data set lists the time that appeared on the game clock for each play. These times begin at 12:00 at the beginning of each period and then count down to 00:00, which marks the end of the period. The first two digits refer to the number of minutes left in the period. The second two digits refer to the number of seconds.

The times have not been parsed as date-time data to R. It would be difficult to record the time data as a date-time object because the data is incomplete: a minutes and seconds element are not sufficient to identify a unique date-time. However, we can store the minutes and seconds information as a *period* object, as defined in Section 5.4, using the `ms()` parse function.

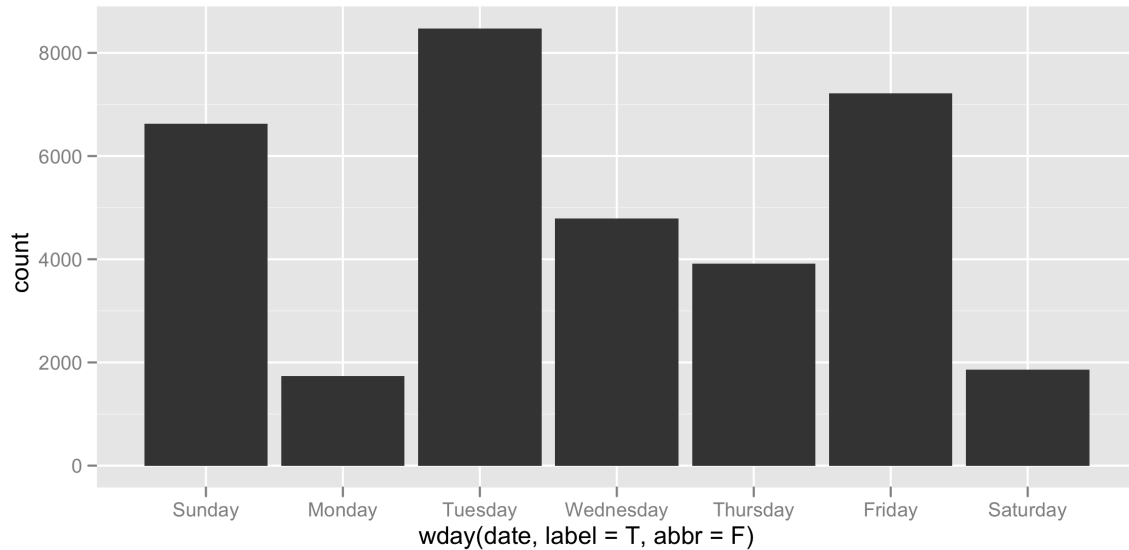


Figure 2: The number of games played per weekday varies.

```
R> lakers$time <- ms(lakers$time)
```

Since periods have relative lengths, it is dangerous to compare them to each other. So we should next convert our periods to *durations*, which have exact lengths.

```
R> lakers$time <- as.duration(lakers$time)
```

This allows us to directly compare different durations. It would also allow us to determine exactly when each play occurred by adding the duration to the *instant* the game began. (Unfortunately, the starting time for each game is not available in the data set). However, we can still calculate when in each game each play occurred. Each period of play is 12 minutes long. At the start of each period, the game clock begins counting down from 12:00. So to calculate how much play time elapses before each play, we subtract the time that appears on the game clock from a duration of 12, 24, 36, or 48 minutes (depending on the period of play). This creates a new duration that records exactly how far into the game each play occurred.

```
R> lakers$time <- eminutes(12) * lakers$period - lakers$time
```

One complication is that some games went into overtime. These games have a fifth period that is only five minutes long. To keep things simple, we will ignore plays that occur in overtime.

```
R> lakers <- subset(lakers, period != 5)
```

ggplot2 does not support plotting *difftimes*, which is the class used by durations. To plot our data, we can extract the integer value of our durations, which will equal the number of

seconds that occurred in each duration (Figure 3).

```
R> qplot(as.integer(time), data = lakers, geom = "histogram", binwidth = 60)
```

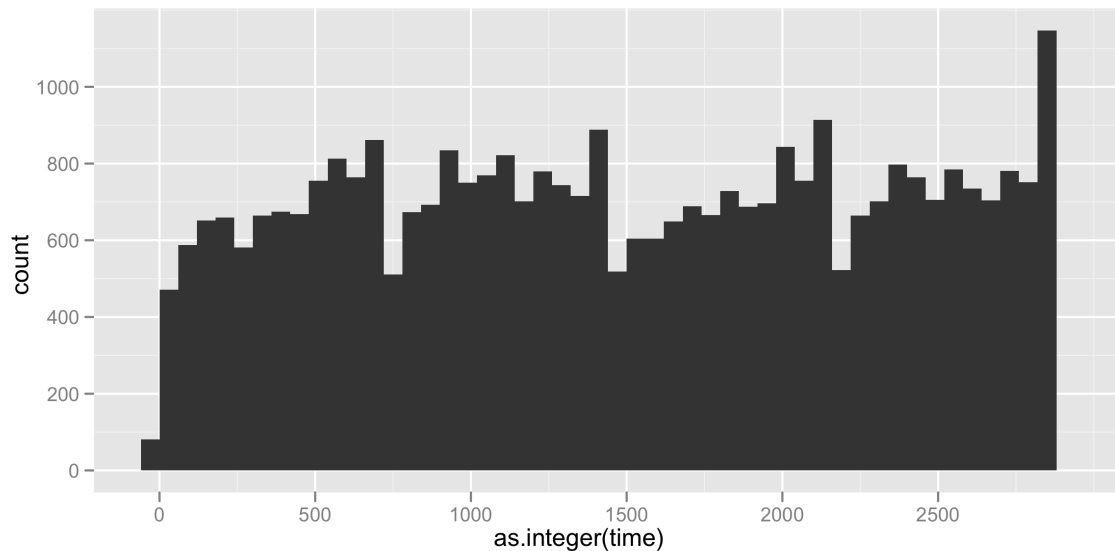


Figure 3: By default, the time of each play is recorded in seconds on the x axis.

Or we can create date-times, which **ggplot2** does support, by adding each of our durations to the same starting instant. This creates a plot whose tick marks are determined by `pretty.date()`. This helper function recognizes the most intuitive binning and labeling of date-time data, which further enhances our graph (Figure 4).

```
R> lakers$minutes <- ymd("2008-01-01") + lakers$time
R> qplot(minutes, data = lakers, geom = "histogram", binwidth = 60)
```

The number of plays peaks within each of the four periods and then plummets at the beginning of the next period, figure 4.

Now lets look more closely at just one basketball game: the first game of the season. This game was played on October 28, 2008. For this game, we can easily model the amounts of time that occurred between each shot attempt.

```
R> game1 <- lakers[lakers$date == ymd("20081028"),]
R> attempts <- game1[game1$type == "shot",]
```

The waiting times between shots will be the time span that occurs between each shot attempt. Since we've recorded the time of each shot attempt as a duration (above), we can record the differences by subtracting the two durations. This automatically creates a new duration whose length is equal to the difference between the first two durations.

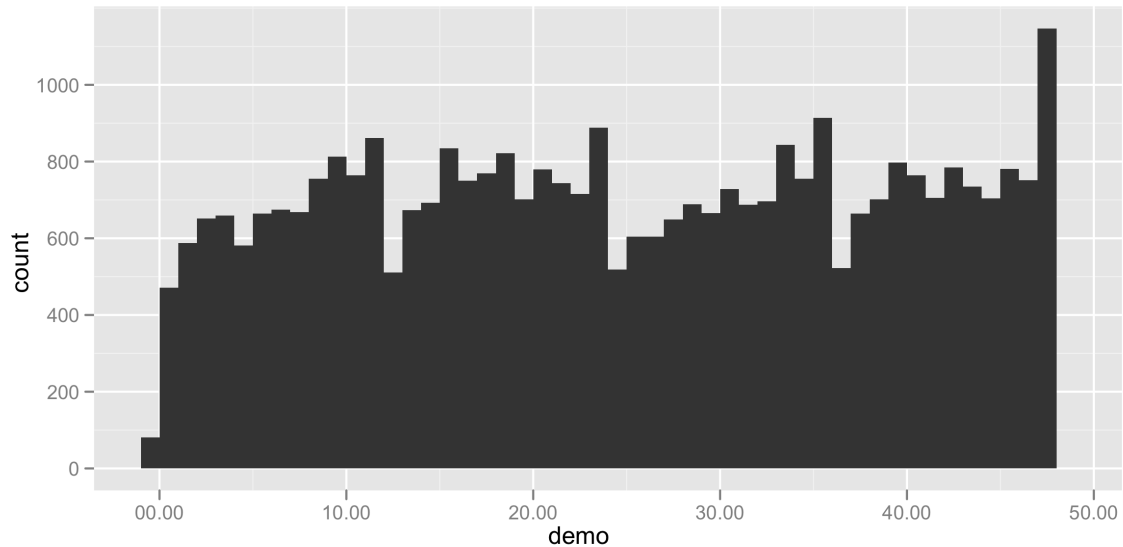


Figure 4: It makes more sense to show the time of each play in minutes.

```
R> attempts$wait <- c(attempts$time[1], diff(attempts$time))
R> qplot(as.integer(wait), data = attempts, geom = "histogram", binwidth = 2)
```

We plot this information in Figure 5. We see that 30 seconds rarely go by without at least one shot attempt, but on occasion up to 60 seconds will pass without an attempt.

We can also examine changes in the score throughout the game. This reveals that the first game of the season was uneventful: the Lakers maintained a lead for the entire game, Figure 6. Note: the necessary calculations are made simpler by the `ddply()` function from the `plyr` package, which `lubridate` automatically loads.

```
R> game1_scores <- ddply(game1, "team", transform, score = cumsum(points))
R> game1_scores <- game1_scores[game1_scores$team != "OFF",]
R> qplot(ymd("2008-01-01") + time, score, data = game1_scores, geom = "line", colour = team)
```

11. Conclusion

Date-times create technical difficulties that other types of data do not. They must be specifically identified as date-time data, which can be difficult due to the overabundance of date-time classes. It can also be difficult to access and manipulate the individual pieces of data contained within a date-time. Math with date-times is often appropriate, but must follow different rules than math with ordinal numbers. Finally, date-related conventions such as daylight savings time and time zones make it difficult to compare and recognize different moments of time.

Base R handles many of these difficulties, but not all. Moreover, base R's date-time methods are often complicated, confusing, and prone to change for different date-time classes.

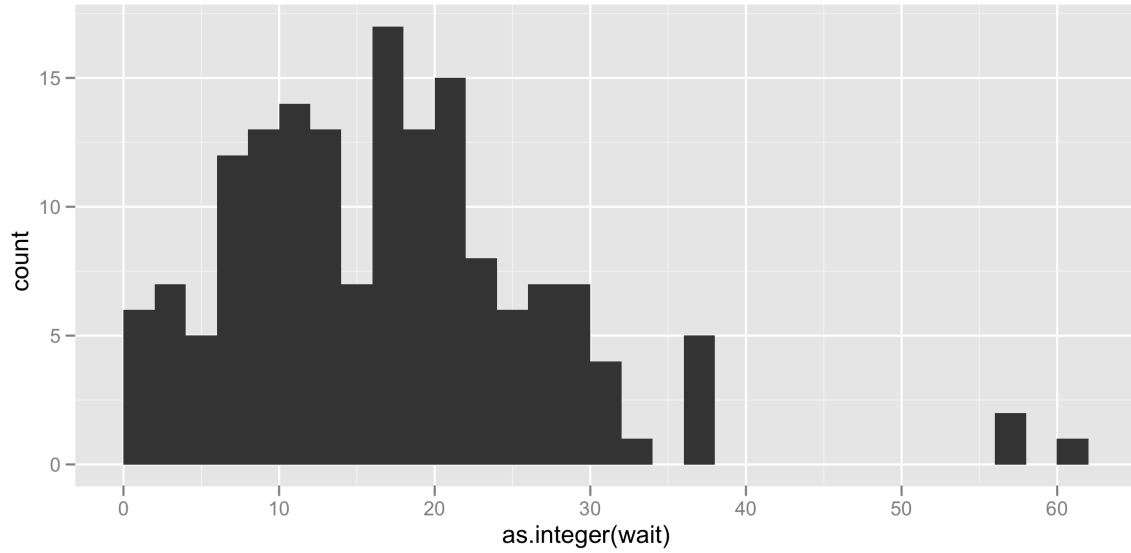


Figure 5: Distribution of wait times between shot attempts.

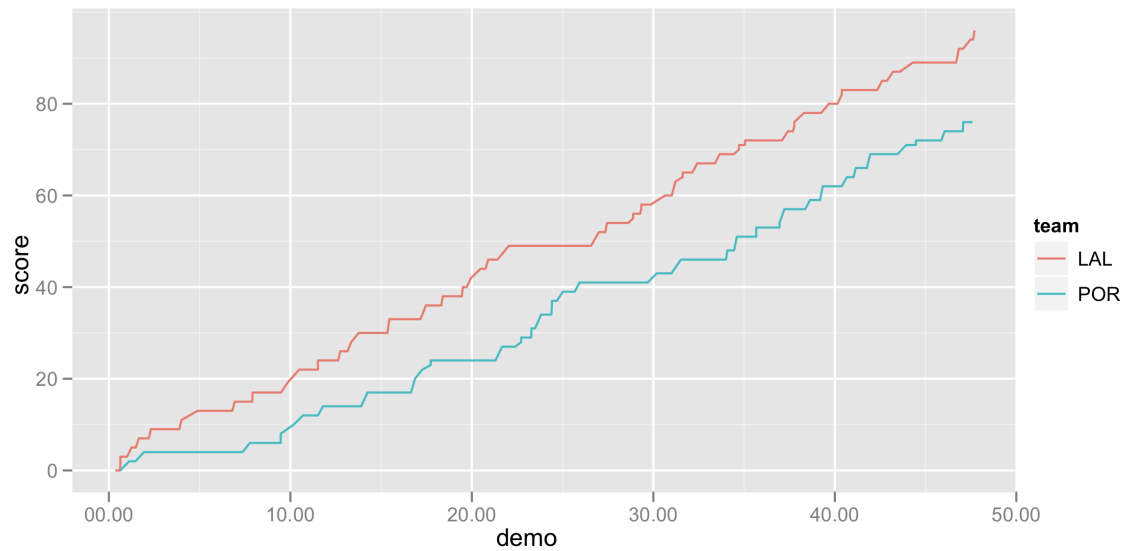


Figure 6: Scores over time for first game of season.

lubridate makes it easier to work with date-time data in R. The package provides a set of standard methods for most common date-time classes. These methods make it simple to parse, manipulate, and perform math on date-time objects. By implementing the time concepts pioneered by the Java based JODA time project, **lubridate** helps researchers perform precise calculations as well as model tricky time-related processes. **lubridate** also makes it simple to switch date-times between time zones and to advantageously use or ignore daylight savings time.

Future work on the **lubridate** package will develop methods for handling partial dates and for modeling re-occurring events, such as stock market opening times, business days, or street cleaning hours. In particular, we hope to create methods for R that work with reoccurring temporal date patterns, which were introduced by Martin Fowler (Fowler 1997) and have been implemented in Ruby by the runt project (Lipper 2008).

Acknowledgements

We'd like to thank the National Science Foundation. This work was supported by the NSF VIGRE Grant, number DMS-0739420.

References

- Armstrong W (2009). *fts: R Interface to tslib (A Time Series Library in C++)*. R package version 0.7.6, URL <http://CRAN.R-project.org/package=fts>.
- Colebourne S, O'Neill BS (2006). "Joda-Time." URL <http://joda-time.sourceforge.net/>.
- Fowler M (1997). "Recurring events for calendars." *n. pub., n. pag*, pp. 0–201.
- Grothendieck G, Petzoldt T (2004). "Date and Time Classes in R." *R News*, 4(1), 32.
- Hallman J (2010). *tis: Time Indexes and Time Indexed Series*. R package version 1.12, URL <http://CRAN.R-project.org/package=tis>.
- James D, Hornik K (2010). *chron: Chronological Objects which Can Handle Dates and Times*. R package version 2.3-35. S original by David James, R port by Kurt Hornik., URL <http://CRAN.R-project.org/package=chron>.
- Lipper M (2008). "runt." URL <http://runt.rubyforge.org/>.
- Parker RJ (2010). "Basketball Geek." URL <http://www.basketballgeek.com/data/>.
- Portfolio and Risk Advisory Group, Commerzbank Securities (2009). *its: Irregular Time Series*. R package version 1.1.8, URL <http://CRAN.R-project.org/package=its>.
- R Development Core Team (2010). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org>.
- Ryan JA, Ulrich JM (2010). *xts: Extensible Time Series*. R package version 0.7-1, URL <http://CRAN.R-project.org/package=xts>.

- Trapletti A, Hornik K (2009). *tseries: Time Series Analysis and Computational Finance*. R package version 0.10-22., URL <http://CRAN.R-project.org/package=tseries>.
- Wickham H (2009). *ggplot2: Elegant Graphics for Data Analysis*. Springer New York. ISBN 978-0-387-98140-6. URL <http://had.co.nz/ggplot2/book>.
- Wuertz D, Chalabi Y (2010). *timeSeries: Rmetrics - Financial Time Series Objects*. R package version 2110.87, URL <http://CRAN.R-project.org/package=timeSeries>.
- Wuertz D, with contributions from Martin Maechler YC, Byers JW, *et al.* (2010). *timeDate: Rmetrics - Chronological and Calendarical Objects*. R package version 2110.88, URL <http://CRAN.R-project.org/package=timeDate>.
- Zeileis A, Grothendieck G (2005). “**zoo**: S3 Infrastructure for Regular and Irregular Time Series.” *Journal of Statistical Software*, **14**(6), 1–27. URL <http://www.jstatsoft.org/v14/i06/>.

Affiliation:

Garrett Grolemond
Rice University
Houston, TX 77251-1892, United States of America
E-mail: grolemond@rice.edu