

Learning *R*

Carl James Schwarz

StatMathComp Consulting by Schwarz
cschwarz.stat.sfu.ca @ gmail.com

Split-Apply-Combine Paradigm
Advanced usage of the *plyr* package

1. Split-apply-combine - advanced *plyr* package
 - 1.1 Passing a function to *ddply*
 - 1.2 Passing a function to *ddply* with additional arguments
 - 1.3 *dlply* and *ldply*
 - 1.4 Parallelization
 - 1.5 Array operations

Split - Apply - Combine

Performing the same analysis
to multiple chunks of your data

Advanced usage of *plyr* package.

Recall

```
res <- plyr::ddply(cereals, "Shelf", plyr::summarize,  
  mean.fat=mean(fat, na.rm=TRUE)  
  ...  
)
```

Passing a function to *ddply()* function

- Sometimes computations are too complex to use *ply::summarize*
- The chunk is passed as DATA.FRAME and traditionally called *x* and is the first argument of the function
- The function can be defined separately or as part of the call (a.k.a. anonymous function)

Split - Apply - Combine - *ddply()* + your own function

Here is a simple function - what does it do?

```
1 mysummary <- function(x){
2   # compute the mean fat and calories and their ratio
3   mean.fat = mean(x$fat, na.rm=TRUE)
4   mean.calories=mean(x$calories, na.rm=TRUE)
5   ratio = mean.calories / mean.fat
6   data.frame(mean.fat, mean.calories, ratio, stringsAsFactors=FALSE)
7 }
```

We test the function:

```
> mysummary(cereal)
  mean.fat mean.calories    ratio
1 1.012987    105.0649 103.7179
```

As always, functions should be *self-contained* and seldom refer to variable not passed as arguments or in the calling environment!

Split - Apply - Combine - *ddply()* + your own function

We pass the function to *plyr::ddply()*

```
1 report <- plyr::ddply(cereal, "shelf", mysummary)
2 report
```

We get a separate summary for each shelf

```
> report <- plyr::ddply(cereal, "shelf", mysummary)
> report
```

	shelf	mean.fat	mean.calories	ratio
1	1	0.60	100.5000	167.50000
2	2	1.00	107.6190	107.61905
3	3	1.25	106.1111	84.88889

As always, functions should be *self-contained* and seldom refer to variable not passed as arguments or in the calling environment!

Split - Apply - Combine - *ddply()* + your own function

Rather than cluttering up the environment with functions that are only used once, it is quite common to use anonymous functions. Simply replace the function name above by the actual function definition:

```
plyr::ddply(dataframe, byvars, function(x)
  { # don't forget the opening brace
  func definition
  res1 <- ..
  res2 <- ...
  res <- data.frame(res1,res2,
                    stringsAsFactors=FALSE)
  return(res)
} # don't forget the closing brace
) # don't forget the closing )
```


Split - Apply - Combine - *ddply()* + your own function

Rather than cluttering up the environment with functions that are only used once, it is quite common to use it anonymous functions. Simply replace the function name above by the actual function definition:

```
1 report <- plyr::ddply(cereal, "shelf", function(x){
2   # compute the mean fat and calories and their ratio
3   mean.fat = mean(x$fat, na.rm=TRUE)
4   mean.calories=mean(x$calories, na.rm=TRUE)
5   ratio = mean.calories / mean.fat
6   data.frame(mean.fat, mean.calories, ratio, stringsAsFactors=FALSE)
7 })
```

We get the same results. Be careful to match up braces and parentheses and don't forget to refer to the chunk as *x*.

Split - Apply - Combine - *ddply()* + your own function - Exercise

Fit a separate regression line for each shelf and report the

- intercept
- slope
- RMSE available from *summary(fit)\$sigma*

Use an external function definition and an anonymous function

Split - Apply - Combine - *ddply()* + your own function - Exercise

```
1 library(plyr)
2 sumstats <- plyr::ddply(cereal, "shelf", function(x) {
3     result <- lm(calories ~ fat, data=x)    # notice use
4     intercept <- coef(result)[1]
5     slope <- coef(result)[2]
6     sigma <- summary(result)$sigma
7     res <- data.frame(intercept, slope, rmse,
8                       stringsAsFactors=FALSE)
9     return(res)
10 })
11 sumstats
```

```
> sumstats
```

	shelf	intercept	slope	rmse
1	1	100.27778	0.3703704	11.76983
2	2	96.78571	10.8333333	9.09777
3	3	91.36752	11.7948718	25.82378

Refer back to the accidents dataset. For each day, compute

- Number of accidents
- Proportion of fatalities
- MEAN weather severity (*Weather_Conditions*). Not really valid but a close approximation)
- Day of the week (0=Sunday)

Use `plyr::summarize` and write your own (anonymous) function.

Plot number of accident over the year with the SIZE of point related to mean weather conditions.

Add loess curve.

```
1 accidents <- read.csv(file.path(... , 'road-accidents-2010.c
2           as.is=TRUE, strip.white=TRUE)
3 # Convert date to internal date format
4 accidents$mydate <- as.Date(accidents$Date,
5           format="%d/%m/%Y")
6 # Create the fatality variable
7 accidents$Fatality <- accidents$Accident_Severity == 1
```

Using *ddply()* and *summarize()*

```
1 naccidents <- plyr::ddply(accidents, "mydate", plyr::summarize(  
2     freq=length(mydate),  
3     pfatal=mean(Fatality),  
4     mean.weather=mean(Weather_Conditions),  
5     dow=format(mydate, "%w")[1])  
6 naccidents[1:5,]
```

```
> naccidents[1:5,]
```

	mydate	freq	pfatal	mean.weather	dow
1	2010-01-01	282	0.014184397	2.262411	5
2	2010-01-02	293	0.030716724	2.740614	6
3	2010-01-03	273	0.014652015	2.857143	0
4	2010-01-04	401	0.002493766	2.518703	1
5	2010-01-05	379	0.002638522	2.936675	2

Split - Apply - Combine - Exercise VI

Using *ddply()* and an explicit function

```
1 naccidents <- plyr::ddply(accidents, "mydate", function(x){
2     freq <- nrow(x)
3     mean.weather <- mean(x$Weather_Conditions)
4     pfatal <- mean(x$Fatality)
5     dow=format(x$mydate, "%w")[1]
6     res <- data.frame(freq, mean.weather, pfatal,
7                       dow, stringsAsFactors=FALSE)
8     return(res)
9 })
10 naccidents[1:5,]
```

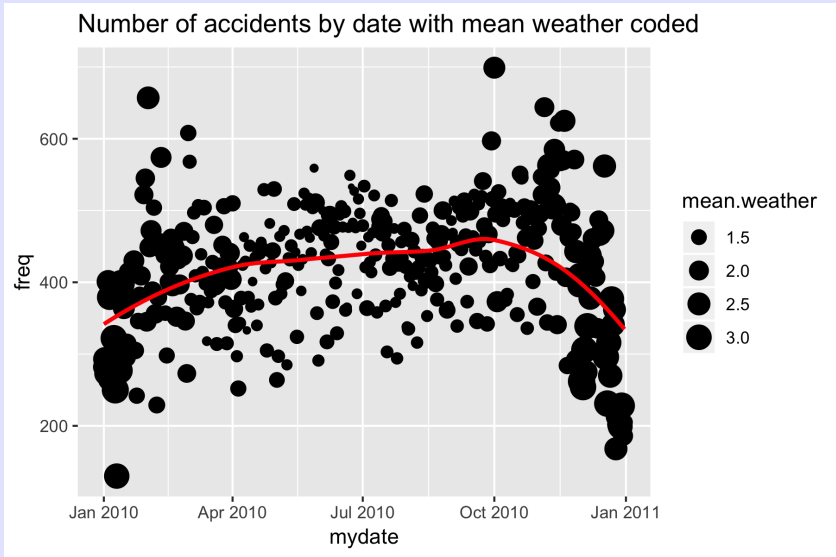
```
> naccidents[1:5,]
```

	mydate	freq	mean.weather	pfatal	dow
1	2010-01-01	282	2.262411	0.014184397	5
2	2010-01-02	293	2.740614	0.030716724	6
3	2010-01-03	273	2.857143	0.014652015	0
4	2010-01-04	401	2.518703	0.002493766	1
5	2010-01-05	379	2.936675	0.002638522	2

Make the plots

```
1 newplot <- ggplot(data=naccidents,  
2                   aes(x=mydate, y=freq ))+  
3   ggtitle("Number of accidents by date with mean weather coo  
4   geom_point( aes(size=mean.weather))+  
5   geom_smooth(method="loess", color="red", se=FALSE)  
6 newplot
```

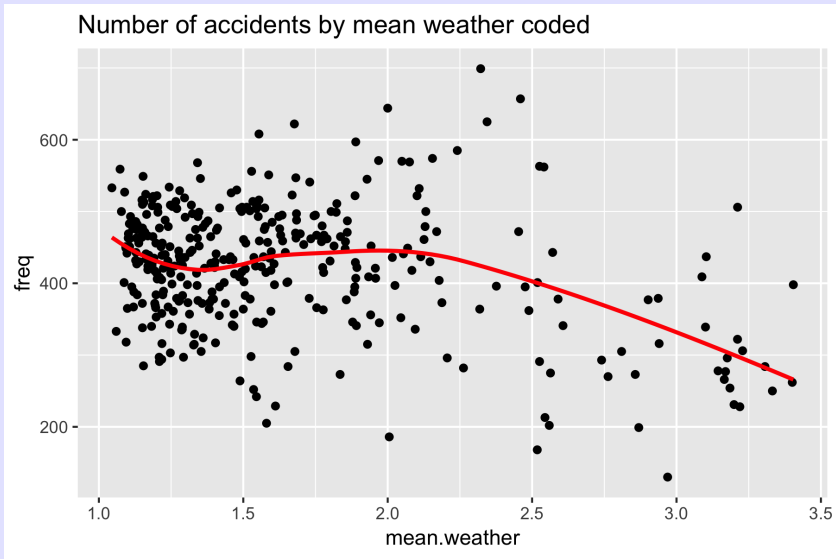

Split - Apply - Combine - Exercise VI



Plot number of accidents vs. mean weather conditions;
Add loess curve

Plot number of accidents vs. mean weather conditions;

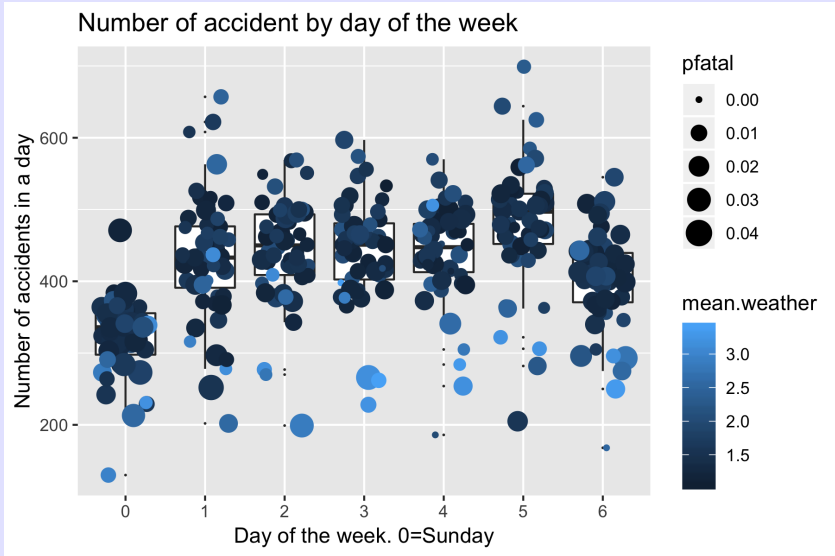
```
1 newplot <- ggplot(data=naccidents,  
2                   aes(x=mean.weather, y=freq))+  
3   geom_point( ) +  
4   geom_smooth(method="loess", color="red", se=FALSE) +  
5   ggtitle("Number of accidents by mean weather coded")  
6 newplot
```



Accident data.

Make a box-plot of number of accident by day of the week coded using proportion of fatalities by the size of the symbol and the mean weather condition by a color gradient.

```
1 newplot <- ggplot(data=naccidents, aes(x=dow, y=freq))+
2   geom_boxplot( ) +
3   geom_jitter(aes(size=pfatal, color=mean.weather),
4     position=position_jitter(w=.3, h=.0))+
5   ggtitle("Number of accident by day of the week")+
6   xlab("Day of the week. 0=Sunday") +
7   ylab("Number of accidents in a day")
8 newplot
```



VERY COMMON PARADIGM IN *R*.

- Virtually unnecessary to use *for* loops in *R* if computations for each chunk are independent and do not depend on other chunks.
- Makes it easy to parallelize your work (routines are set up to use multiple machines)
- Most common usage is *ddply()*

Split - Apply - Combine - Summary (Simple)

Most simple usage is with *ddply()* and *summarize()*,
Sometimes it is more convenient write your own function

```
1 sumstat <- plyr::ddply( dataframe, "chunking variable",
2                       function(x){
3                           res1 <- function of x$...
4                           res2 <- function of x$...
5                           res <- data.frame(res1, res2, ...,
6                                             stringsAsFactors=FALSE)
7                           return(res)
8                       })
```

Passing additional arguments to *ddply()* functions.
Sometimes you need to pass additional variables other than the chunk to be processed.

Example: Refer back to the cereal dataset. For each shelf group (and for an “arbitrary” variable), compute

- Number of observations
- Number of observations with missing values
- Mean of the variable
- SD of the variable

Split - Apply - Combine - Advanced

```
1 sumstat <- function(x, var){
2   # Compute some summary statistics for a data frame
3   values <- x[,var] # extract the variable values
4   n <- length(values)
5   nmiss <- sum(is.na(values))
6   mean <- mean(values, na.rm=TRUE)
7   sd <- sd(values, na.rm=TRUE)
8   res <- data.frame(n,nmiss,mean,sd,
9                     stringsAsFactors=FALSE)
10  return(res)
11 } # end of sumstat
12
13
14 sumstat(cereal, "calories")
15 sumstat(cereal, "weight")
```

But how are the second (and additional arguments passed to `ddply()`?

```
1 res <- plyr::ddply( dataframe, "chunking variable",
2                     functionname,
3                     y=xxx, z=xxxx)
4
5 res <- plyr::ddply( dataframe, "chunking variable",
6                     function(x, y, z){
7                         res1 <- function of $x$, $y$, and $z
8                         res2 <- function of $x$, $y$, and $z
9                         res <- data.frame(res1, res2, ...,
10                                           stringsAsFactors=FALSE)
11                         return(res)
12                     }, y=xxx, z=xxxx)
13
14 plyr::ddply(cereal, "shelf", sumstat, var="calories")
15 plyr::ddply(cereal, "shelf", sumstat, var="weight")
```

Notice the placement of the variables in the function header and the calling location.

Write a function that takes a data frame and a variable name and

- Find the sample size, number of missing values, mean, its se, and a 95% normal-based confidence interval.
- Bonus - allow for different size of confidence limits, e.g. 90% confidence interval.

Either use the *t.test()* or *lm(y ~ 1)* or code yourself using sample size and t-distribution

Sample output:

```
> my.simple.summary(cereal, "fat")
  variable  n nmiss      mean      sd      se conflevel
1      fat 77     0 1.012987 1.006473 0.1146982      0.95 0
> my.simple.summary(cereal, "fat", conflevel=.90)
  variable  n nmiss      mean      sd      se conflevel
1      fat 77     0 1.012987 1.006473 0.1146982      0.9 0
> plyr::ddply(cereal, "shelf", my.simple.summary, variable="weight")
  shelf variable  n nmiss      mean      sd      se conflevel
1     1     weight 20     0 0.991500 0.03801316 0.00850000
2     2     weight 21     0 1.015714 0.07201190 0.01571429
3     3     weight 36     2 1.062353 0.21450519 0.03678734
> plyr::ddply(cereal, "shelf", my.simple.summary, variable="fat")
  shelf variable  n nmiss mean      sd      se conflevel
1     1      fat 20     0 0.60 0.7539370 0.1685854      0.9
2     2      fat 21     0 1.00 0.7745967 0.1690309      0.9
3     3      fat 36     0 1.25 1.1801937 0.1966989      0.9
```

Using *dply()* and *ldply()*

- It is convenient to do ALL computations for a chunk, return a list, and then extract from the list a needed rather than having several different function.
- Some output (like plots) cannot be stored in data frames.

Example: Refer back to the cereal dataset. For each shelf group (and for two “arbitrary” variable), compute

- Plot of Y vs. X (use *aes_string()* in *ggplot()*)
- Regression of Y on Y . Use *lm(x[, Yvar] ~ x[, Xvar])*
- Return both in a list

```
sumstat(cereal, "calories", "fat")
```

should return a list with 2 elements.

Using *dply()* and *ldply()*

```
1 sumstat <- function(x, Yvar, Xvar){
2   # Do the plot (use aes_string)
3   plot <- ggplot(data=x, aes_string(x=Xvar, y=Yvar))+
4     ggtitle(paste("Scatterplot of ", Yvar, " vs. ", Xvar, ))+
5     geom_point(position=position_jitter(h=.1, w=.1))+
6     geom_smooth(method="lm", se=FALSE)
7   fit <- lm( x[,Yvar] ~ x[, Xvar], data=x)
8   list(plot=plot, fit=fit)
9 }
10
11 res<- sumstat(cereal, "calories", "fat")
12 length(res)
```


Using *dlply()* and *ldply()*

- Now use *dlply()* to make a list of lists, once for each shelf

```
1 res <- plyr::dlply(cereal, "shelf", sumstat,  
2                   Yvar="calories", Xvar="fat")  
3 length(res)  
4 res[[1]]  
5 res[[2]]  
6 res[[3]]
```

Using `dlply()` and `ldply()`

- Now use `ldply()` to make a data frame of slope and se

```

1 se.summary <- plyr::ldply(res, function(x){
2     # x is now the list of the flot and the fit
3     slope <- summary(x$fit)$coefficients[2,]
4     slope
5 })
6 se.summary

```

	shelf	Estimate	Std. Error	t value	Pr(> t)
1	1	0.3703704	3.581444	0.1034137	0.9187781230
2	2	10.8333333	2.626300	4.1249412	0.0005759948
3	3	11.7948718	3.698559	3.1890452	0.0030617682

Using `dlply()` and `ldply()`

- Now use `ldply()` to extract the plots

```
1 plyr::l_ply(res, function(x){
2   plot(x$plot) # needed because within a function
3 })
```

- Notice use of `l_ply()` because no output returned.
- Notice use of `plot()` WITHIN function to force display.
- All plots are sent to the plot window in *Rstudio*.
- It is possible to send the plots to a pdf file (see code).

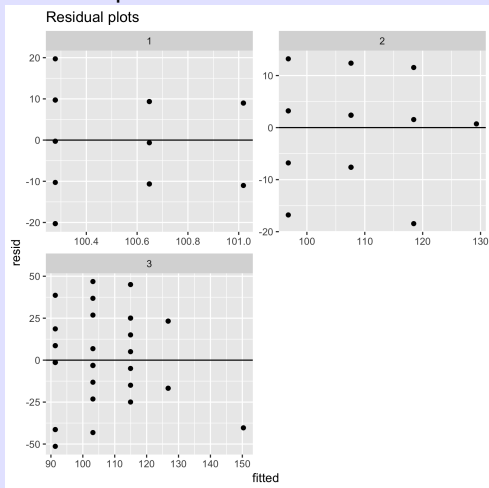
Create residual displays using *dlply()* and *ldply()*

```
1 resid <- ldply(res, function(x){
2   data.frame(resid= resid(x$fit), fitted=fitted(x$fit))
3 })
4 head(resid)
```

	shelf	resid	fitted
1	1	8.981481	101.0185
2	1	-10.648148	100.6481
3	1	8.981481	101.0185

```
1 ggplot(data=resid, aes(x=fitted, y=resid))+  
2   ggtitle("Residual plots")+  
3   geom_point()+  
4   geom_hline(yintercept=0)+  
5   facet_wrap(~shelf, ncol=2, scales="free")
```

Residual plot for each shelf:



All of the routines in the *plyr* package allow parallelization, i.e. using separate cores on your machine.

- Set up cores to be used
- Set `.parallel=TRUE` in the call
- Close the cores used.

Set up the cores

```
1 doParallel <- TRUE # should I set up parallel processing o
2
3 if(doParallel) {
4   library(doMC) # for parallel model fitting
5   # see http://viktoriawagner.weebly.com/blog/five-steps-to
6   detectCores()
7   cl <- makeCluster(4)
8   # Need to export some libraries to the cluster
9   # see http://stackoverflow.com/questions/18981932/logging
10  clusterEvalQ(cl, library(unmarked))
11  registerDoMC(5)
12 }
```


Run the `plyr::function()` in parallel.

```
1 Sys.time()
2 res <- plyr::ddply(cereal, "shelf", plyr::summarize,
3                   mean=mean(shelf[1]*(1:100000000), na.rm=
4                               .parallel=FALSE)
5 Sys.time()
6
7 Sys.time()
8 res <- plyr::ddply(cereal, "shelf", plyr::summarize,
9                   mean=mean(shelf[1]*(1:100000000), na.rm=
10                              .parallel=doParallel)
11 Sys.time()
```

Stop the cores.

```
1 if(doParallel) stopCluster(cl) # stop parallel processing
```

Row or Column operations on a MATRIX or ARRAY (less common)
Note MATRIX differs from a data.frame because all values must have same type.

ARRAY is a 3+ dimensional object.

```
1 mat <- matrix(1:30, nrow=6)
2 mat
```

```
> mat
      [,1] [,2] [,3] [,4] [,5]
[1,]    1    7   13   19   25
[2,]    2    8   14   20   26
[3,]    3    9   15   21   27
[4,]    4   10   16   22   28
[5,]    5   11   17   23   29
[6,]    6   12   18   24   30
```

Split - Apply - Combine - Advanced

Row or Column operations on a MATRIX or ARRAY.

```
1 aapply(mat, 1, sum)
2 aapply(mat, 2, mean)
3 aapply(mat, 1, function(x){
4   res <- prod(sin(x))
5   return(res)
6 })
```

```
> aapply(mat, 1, sum)
```

```
  1  2  3  4  5  6
65 70 75 80 85 90
```

```
> aapply(mat, 2, mean)
```

```
  1    2    3    4    5
3.5  9.5 15.5 21.5 27.5
```

```
> aapply(mat, 1, function(x)....
```

```
                1                2                3                4
-0.0046076865  0.6204102446  0.0302615870  0.0002842306 -0.5
```

Refer back to the accident dataset.

- For each month, compute the number of days, weekdays and weekends. Hint: Use the *unique()* function on the dates within each month. Why do I want all three values?
- For each month, compute the total number of accidents with injury, those on weekends, and those on weekdays. Again, why do I want all three values?
- For each month, find the ratio of the number of accidents on weekday to weekends.
- Plot these over the year
- Add a suitable comparison line if accidents were uniformly spread over the days of the week. Note that the number of weekends and weekdays differs among months.

```
1 ... read accident data ....
2 ... convert dates to internal R format ...
3
4 # get the month for each accident date
5 accidents$month <- as.numeric(format(
6     accidents$mydate, "%m"))
```

Split - Apply - Combine - Advanced - Exercise VIII

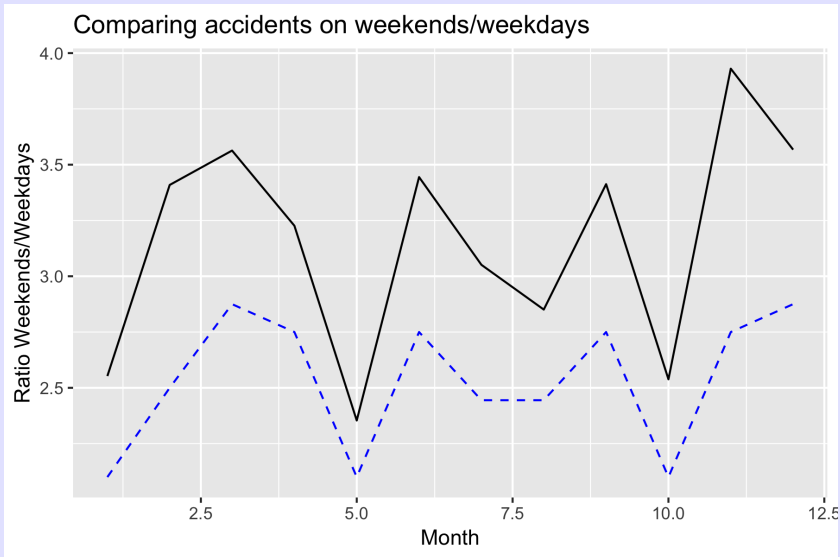
```
1 mysummary <- function(accidents){
2   # Compute the number of weekend and weekdays in the month (
3   # Compute the number of accidents on weekend/weekdays
4   # Report the two ratio.
5   DaysOfMonth <- unique(accidents$mydate)
6   DaysOfWeeks <- format(DaysOfMonth, "%w")
7   nDays      <- length(DaysOfMonth)
8   nWeekdays <- sum(DaysOfWeeks %in% 1:5)
9   nWeekends  <- sum(DaysOfWeeks %in% c(0,6))
10
11  AccDaysOfWeek <- format(accidents$mydate, "%w")
12  nAccTotal     <- length(accidents$mydate)
13  nAccWeekdays <- sum(AccDaysOfWeek %in% 1:5)
14  nAccWeekends  <- sum(AccDaysOfWeek %in% c(0,6))
15
16  rAccWdWe <- nAccWeekdays/nAccWeekends
17  rDaysWdWe <- nWeekdays/nWeekends
18
19  res <- data.frame(nDays,
20                   nWeekdays,
```

```
1 testdata <- subset(accidents, accidents$month == 1)
2 dim(testdata)
3
4 mysummary(testdata)
5
6 > mysummary(testdata)
7       nDays    nWeekdays    nWeekends    nAccTotal  nAccWeel
8       31.000000    21.000000    10.000000  10637.000000    7643.00
```

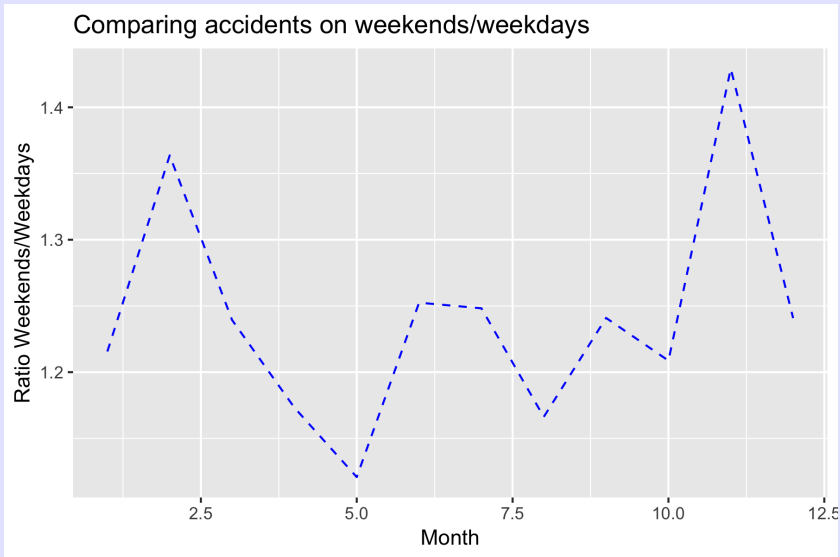

Split - Apply - Combine - Advanced - Exercise VIII

```
1 results <- plyr::ddply(accidents, "month", mysummary)
2 results
3
4      month nDays nWeekdays nWeekends nAccTotal nAccWeekdays n
5 1         1    31         21         10     10637         7643
6 2         2    28         20          8     11724         9065
7 3         3    31         23          8     13165        10280
8 4         4    30         22          8     12248         9350
9 5         5    31         21         10     13220         9278
10 6         6    30         22          8     13644        10574
11 7         7    31         22          9     13527        10188
12 8         8    31         22          9     13027         9644
13 9         9    30         22          8     13904        10753
14 10        10    31         21         10     14429        10351
15 11        11    30         22          8     14544        11594
16 12        12    31         23          8     10345         8080
```

```
1 newplot <- ggplot(data=results, aes(x=month, y=rDaysWdWe))+
2   ggtitle("Comparing accidens on weekends/weekdays")+
3   xlab("Month")+ylab("Ratio Weekends/Weekdays")+
4   geom_line(group=1, color="blue", linetype=2)+
5   geom_line(aes(y=rAccWdWe,group=1))
6 newplot
```



```
1 newplot <- ggplot(data=results, aes(x=month, y=rAccWdWe/rDay
2   ggtitle("Comparing accidens on weekends/weekdays")+
3   xlab("Month")+ylab("Ratio Weekends/Weekdays")+
4   geom_line(group=1, color="blue", linetype=2)
5 newplot
```



General steps for the Split-Apply-Combine paradigm.

- What form is data in? Array, Dataframe, List?
- What form should results be in? Array, Dataframe, List, NULL (for plots)
- Create function to do the application; test with a subset of the data;

```
1 myfunction <- function() {} # define your function
2 testdata <- subset( mydata,
3                   mydata$variable == testvalue)
4 myfunction(testdata)
```

- Run the member of the *plyr* package.