# Learning *R*

Carl James Schwarz

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

## Split-Apply-Combine Paradigm
Introduction to *plyr* package

Split - Apply - Combine
Performing the same analysis
to multiple chunks of your data
*R*'s implementation of PivotTables (and
more) in Excel

## Split - Apply - Combine
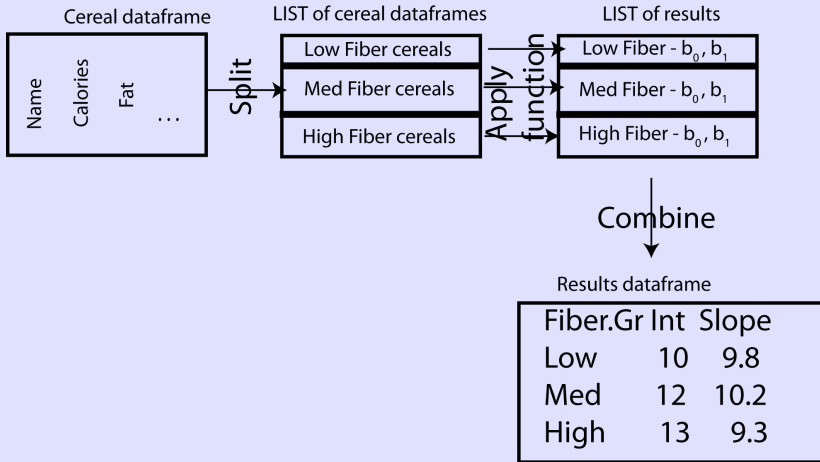
Split-Apply-Combine

- **Split** up a big data frame
- **Apply** a function to each piece
- **Combine** the results together

Examples

- Compute the mean calories/serving for each display shelf.
- Compute number of accidents and p(fatality) for each day in the year.
- Fit a separate regression line to different fiber groups.
- Do a separate analysis for each year of accident data.

Find the slope and intercept of the regression of Calories vs. Fat for each Fiber Group in the cereal dataframe.



Cereal dataframe

Name | Calories | Fat | ...

Split

LIST of cereal dataframes

Low Fiber cereals
Med Fiber cereals
High Fiber cereals

Apply function

LIST of results

Low Fiber - $b_0$, $b_1$
Med Fiber - $b_0$, $b_1$
High Fiber - $b_0$, $b_1$

Combine

Results dataframe

| Fiber.Gr | Int | Slope |
|----------|-----|-------|
| Low      | 10  | 9.8   |
| Med      | 12  | 10.2  |
| High     | 13  | 9.3   |

## Split - Apply - Combine

Base *R* procedures (AVOID)

- *by()* takes data.frame → list
- *split()* takes data.frame → list
- *lapply()* takes list → list
- *sapply()* takes list → vector or matrix

*plyr* package (much more logically arranged) (RECOMMENDED)

- *xyply()* where *x* and *y* are d=data.frame, l=list, a=array, _=nothing
- Hadley Wickham (2011).
  The Split-Apply-Combine Strategy for Data Analysis.
  Journal of Statistical Software, 40(1), 1-29.
  http://www.jstatsoft.org/v40/i01/.
- *dplyr* package - more advanced and only for data.frames.

## Split - Apply - Combine

AVOID: Base $R$ (input data structure (left); output data structure (top))

|  | array | data frame | list | nothing |
|---|---|---|---|---|
| array | apply |  |  |  |
| data frame |  | *aggregate* | by |  |
| list | sapply |  | lapply |  |
| n replicates | replicate |  | replicate |  |
| function arguments | mapply |  | mapply |  |

## Split - Apply - Combine

USE: *plyr* package (input data structure (left); output data structure (top))

|                      | array  | data frame | list   | nothing |
|----------------------|--------|------------|--------|---------|
| array                | aaply  | adply      | alply  | a_ply   |
| data frame           | daply  | **ddply**  | **dlply** | d_ply   |
| list                 | laply  | **ldply**  | llply  | l_ply   |
| n replicates         | raply  | rdply      | **rlply** | r_ply   |
| function arguments   | maply  | mdply      | **mlply** | m_ply   |

## Split - Apply - Combine - *ddply()* + *summarize()*

Cereal dataset.
Find the number of cereals and the mean calories/serving for each
shelf

```
1  cereal <- read.csv(file.path(..., 'cereal.csv'),
2                  header=TRUE, as.is=TRUE,
3                  strip.white=TRUE)
4  cereal[1:5,]
5
6  library(plyr)
7  sumstats <- plyr::ddply(cereal, "shelf", plyr::summarize,
8              ncereal=length(name),
9              mean.calories=mean(calories))
10 sumstats
```

```
> sumstats
  shelf ncereal mean.calories
1   1      20      100.5000
2   2      21      107.6190
3   3      36      106.1111
```

**CAUTION:** Because of conflicts between the *plyr* and *dplyr* packages, ALWAYS

- *plyr::* before the function name
- *plyr::summarize* - this is particularly important.

# Split - Apply - Combine - *ddply() + summarize()* Exercise I

Find the following quantities for each shelf:

- Standard deviation of calories/serving
- Mean number of calories from fat (1 g of fat has 9 calories)
- Mean proportion of calories from fat of total calories.
- Mean weight/serving

```
> sumstats
  shelf std.calores mean.fcal mean.pcal.fat  mean.wt
1     1    11.45931      5.40    0.05404545 0.991500
2     2    12.20851      9.00    0.07986014 1.015714
3     3    29.01012     11.25    0.09859932       NA
```

```
1  library(plyr)
2  sumstats <- plyr::ddply(cereal, "shelf", plyr::summarize,
3              std.calores=sd(calories),
4              mean.fcal = mean(fat*9),
5              mean.pcal.fat = mean( fat*9 / calories),
6              mean.wt=mean(weight))
7  sumstats

   > sumstats
     shelf std.calores mean.fcal mean.pcal.fat  mean.wt
   1    1     11.45931      5.40    0.05404545 0.991500
   2    2     12.20851      9.00    0.07986014 1.015714
   3    3     29.01012     11.25    0.09859932       NA
```

Revise to account for missing values:

```
> sumstats
  shelf std.calores mean.fcal mean.pcal.fat  mean.wt
1     1    11.45931      5.40    0.05404545 0.991500
2     2    12.20851      9.00    0.07986014 1.015714
3     3    29.01012     11.25    0.09859932 1.062353
```

Revise to account for missing values:

```
1  library(plyr)
2  sumstats <- plyr::ddply(cereal, "shelf", plyr::summarize,
3                std.calores=sd(calories),
4                mean.fcal = mean(fat*9),
5                mean.pcal.fat = mean( fat*9 / calories),
6                mean.wt=mean(weight, na.rm=TRUE))
7  sumstats
```

```
> sumstats
  shelf std.calores mean.fcal mean.pcal.fat mean.wt
1    1    11.45931     5.40    0.05404545 0.991500
2    2    12.20851     9.00    0.07986014 1.015714
3    3    29.01012    11.25    0.09859932 1.062353
```

Fit a separate regression line between calories and fat and report the intercept and slope for each shelf.

Recall line for ALL of data is found as:

```
result <- lm(calories ~ fat, data=cereal)
summary(result)
coef(result)
coef(result)[1]
coef(result)[2]

> sumstats
  shelf intercept      slope
1     1 100.27778  0.3703704
2     2  96.78571 10.8333333
3     3  91.36752 11.7948718
```

```
1  library(plyr)
2  sumstats <- plyr::ddply(cereal, "shelf", plyr::summarize,
3             intercept=coef(lm(calories ~fat))[1],
4             slope    =coef(lm(calories ~fat))[2])
5  sumstats
```

```
> sumstats
  shelf intercept      slope
1     1 100.27778  0.3703704
2     2  96.78571 10.8333333
3     3  91.36752 11.7948718
```

A better method will be demonstrated later that doesn't require repeated model fitting.

# Split - Apply - Combine - *ddply()* + *summarize()* Exercise III

Fit a separate regression line between calories and fat and report the intercept and slope for each shelf.

Recall line for ALL of data is found as:

```
result <- lm(calories ~ fat, data=cereal)
summary(result)
coef(result)
coef(result)[1]
coef(result)[2]

> sumstats
  shelf intercept      slope
1     1 100.27778  0.3703704
2     2  96.78571 10.8333333
3     3  91.36752 11.7948718
```

Refer to *road-accidents-2010.csv* file in SampleData.

- Read data into *R*.
- Convert input date to internal *R* dates.
- Find number of accidents by day of year (use *ddply()* and *summarize()* in *plyr* package)
- Plot # accidents/day by day of year.
- Fit a *lowess()* smoother to data using *geom_smooth()*

Look at number of accident by day of the week

- Extract day of the week using *format()* or *weekdays()* functions.
- Use *geom_boxplot()* as seen earlier

```
1  # The accident data
2  accidents <- read.csv(file.path(... ,'road-accidents-2010.cs
3              header=TRUE,
4              as.is=TRUE, strip.white=TRUE)
5  accidents[1:5,]
6  str(accidents)

  > accidents[1:5,]
  .....
    Accident_Severity Number_of_Vehicles Number_of_Casualties
  1                 3                  2                    1
  2                 3                  1                    1

  > str(accidents)
  'data.frame': 154414 obs. of  33 variables:
  ...
   $ Date   : chr  "11/01/2010" "11/01/2010" "12/01/2010" "02/
  ...
```

## Split - Apply - Combine - *ddply()* + *summarize()* Exercise IV

```
1
2  # Convert date to internal date format
3  accidents$mydate <- as.Date(accidents$Date,
4            format="%d/%m/%Y")
5  sum(is.na(accidents$mydate))
6  accidents[1:5,]
7  str(accidents)

   > accidents[1:5,]
   ...
     Urban_or_Rural_Area Did_Police_Officer_Attend_Scene_of_Acc
   1                   1
   2                   1
   > str(accidents)
   'data.frame': 154414 obs. of  33 variables:
    $ Date     : chr "11/01/2010" "11/01/2010" "12/01/2010"
    $ mydate   : Date, format: "2010-01-11" "2010-01-11" "2010
   >
```

```
1  # Summarize number of accidents by date
2  library(plyr)
3  naccidents <- plyr::ddply(accidents, "mydate",
4                            plyr::summarize,
5                            freq=length(Accident_Index))
6  naccidents[1:5,]
7  str(naccidents)

   > naccidents[1:5,]
        mydate freq
   1 2010-01-01  282
   2 2010-01-02  293
   ...
   > str(naccidents)
   'data.frame': 365 obs. of  2 variables:
    $ mydate: Date, format: "2010-01-01" "2010-01-02" "2010-01-
    $ freq  : int  282 293 273 401 379 266 284 322 250 130 ...
```
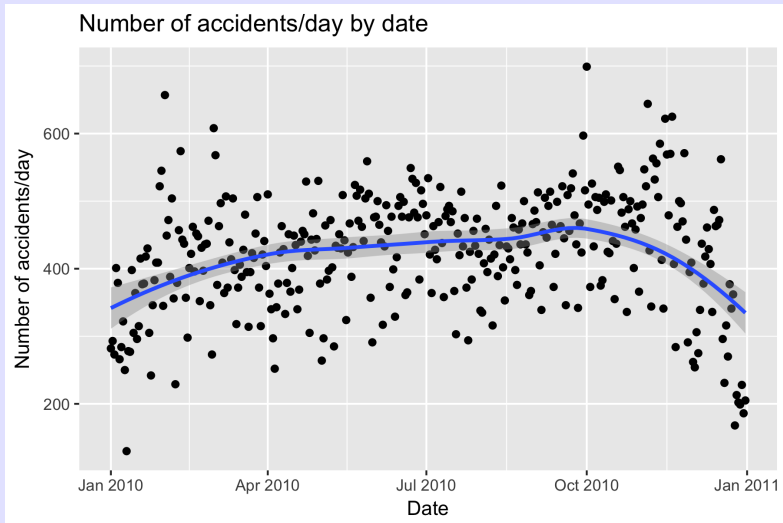
```
1  plotnacc <- ggplot(data=naccidents, aes(x=mydate, y=freq))+
2      ggtitle("Number of accidents/day by date")+
3      xlab("Date")+ylab("Number of accidents/day")+
4      geom_point()+
5      geom_smooth()
6  plotnacc
```

Number of accidents/day by date

Refer to *road-accidents-2010.csv* file in SampleData.

- Create 0/1 variable if fatality occurs (no or yes; check codebook for *Accident_Severity*).
  Use the magic incantation of *recode()* function in *car* package.
- Find proportion of accidents with fatality by day of year
  - The mean of a 0/1 variable is the proportion.
    Use the magic incantation of *ddply()* and *summarize()* in the *plyr* package.
- Plot proportion of fatalities by day of year.
- Fit a *lowess()* smoother to data from *geom_smooth()*
- Plot proportion of fatalities by day of the week
  - Hint: Extract weekday using *format()*.
  - Hint: Use *geom_boxplot()* as seen earlier with some jittering and notches.

```
1  names(accidents)
2  unique(accidents$Accident_Severity)
3  library(car)
4  accidents$Fatality <- recode(accidents$Accident_Severity,
5                 ' 1=1; 2:hi=0')
6  accidents[1:5, c("Accident_Severity", "Fatality")]
7  xtabs(~Fatality + Accident_Severity, data=accidents)

   > accidents[1:5, c("Accident_Severity", "Fatality")]
     Accident_Severity Fatality
   1                 3        0
   2                 3        0

   > xtabs(~Fatality + Accident_Severity, data=accidents)
           Accident_Severity
   Fatality     1      2        3
         0      0  20440   132243
         1   1731      0        0
```

## Split - Apply - Combine - *ddply()* + *summarize()* Exercise V

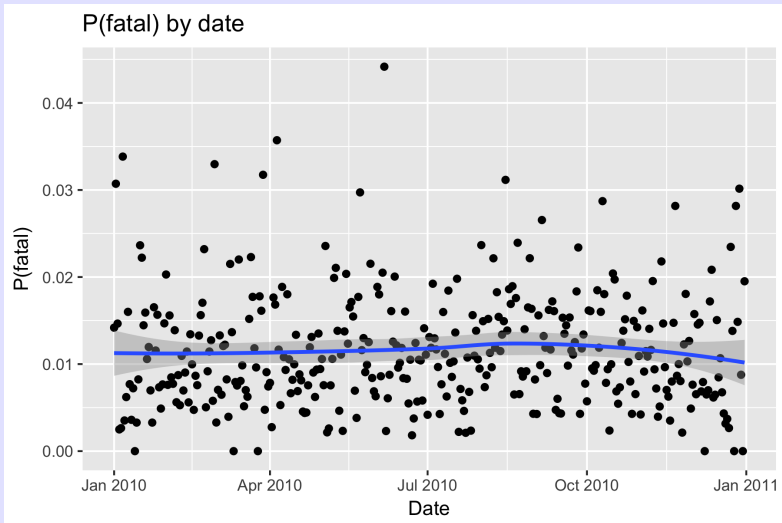The *summarize()* and *ddply()* functions in *plyr* package are quite useful for simple summaries by groups.
Example of the Split-Apply-Combine paradigm to be explained later.

```
1 library(plyr)
2 pfatal.df <- plyr::ddply(accidents, "mydate", plyr::summariz
3              freq=length(mydate),
4              pfatal=mean(Fatality))
5 pfatal.df[1:5,]
```

```
> pfatal.df[1:5,]
       mydate freq      pfatal
1 2010-01-01  282 0.014184397
2 2010-01-02  293 0.030716724
3 2010-01-03  273 0.014652015
4 2010-01-04  401 0.002493766
5 2010-01-05  379 0.002638522
```
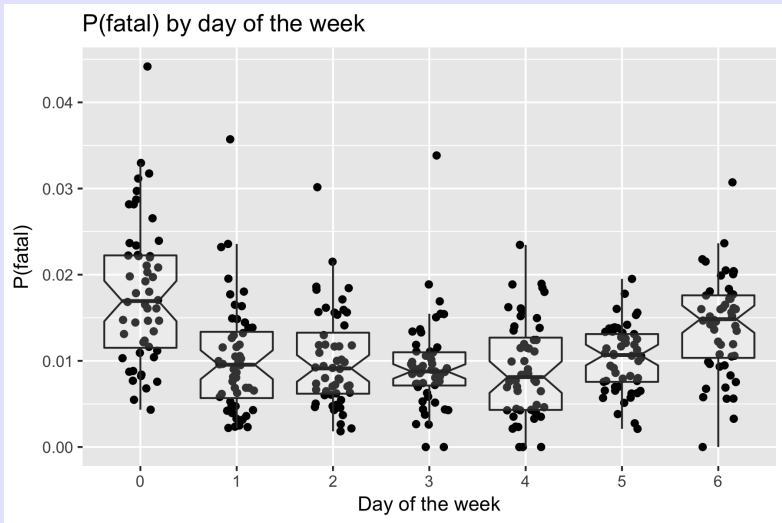
```
1  plotpfatal <- ggplot(data=pfatal.df,
2                  aes(x=mydate, y=pfatal))+
3       ggtitle("P(fatal) by date")+
4       xlab("Date")+ylab("P(fatal)")+
5       geom_point()+
6       geom_smooth()
7  plotpfatal
```

```
 1  # Extract day of the week - leave as character
 2  pfatal.df$weekday <- format(pfatal$mydate, format="%w") # l
 3  pfatal.df[1:10,]
 4
 5  plotpfatal2 <- ggplot(data=pfatal.df, aes(x=weekday, y=pfata
 6    ggtitle("P(fatal) by day of the week")+
 7    xlab("Day of the week")+ylab("P(fatal)")+
 8    geom_point(position=position_jitter(w=0.2))+
 9    geom_boxplot(notch=TRUE, alpha=0.2)
10  plotpfatal2
```

P(fatal) by day of the week

## Split - Apply - Combine - Exercise

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*

Plot number of accidents 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.csv
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
```

# Split - Apply - Combine - Exercise VI

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

```
1  naccidents <- plyr::ddply(accidents, "mydate", plyr::summar:
2                 freq=length(mydate),
3                 pfatal=mean(Fatality),
4                 mean.weather=mean(Weather_Conditions),
5                 dow=format(mydate, "%w")[1])
6  naccidents[1:5,]
```
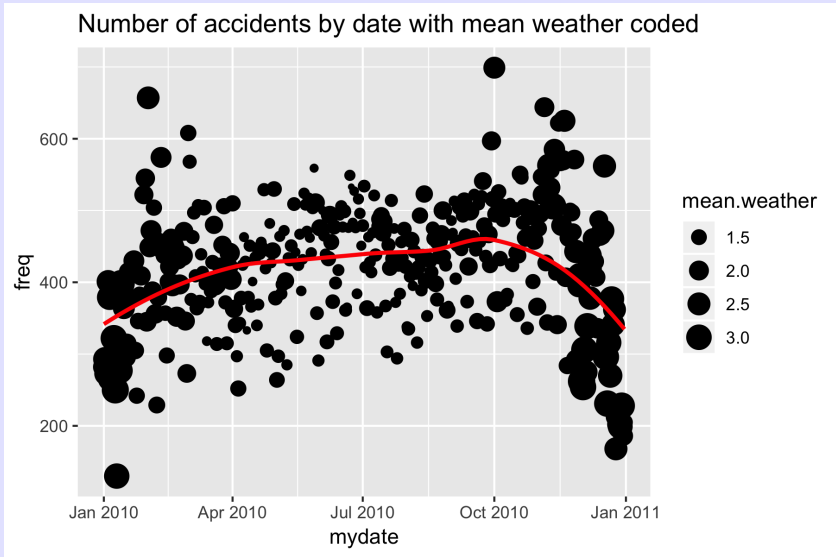
```
> naccidents[1:5,]
      mydate freq      pfatal mean.weather dow
1 2010-01-01  282 0.014184397     2.262411   5
2 2010-01-02  293 0.030716724     2.740614   6
3 2010-01-03  273 0.014652015     2.857143   0
4 2010-01-04  401 0.002493766     2.518703   1
5 2010-01-05  379 0.002638522     2.936675   2
```

Make the plots

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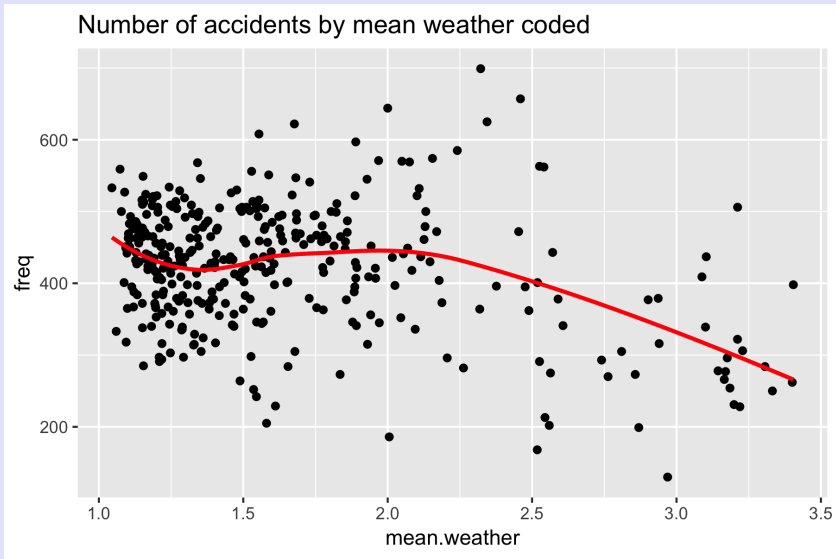
Number of accidents by date with mean weather coded

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

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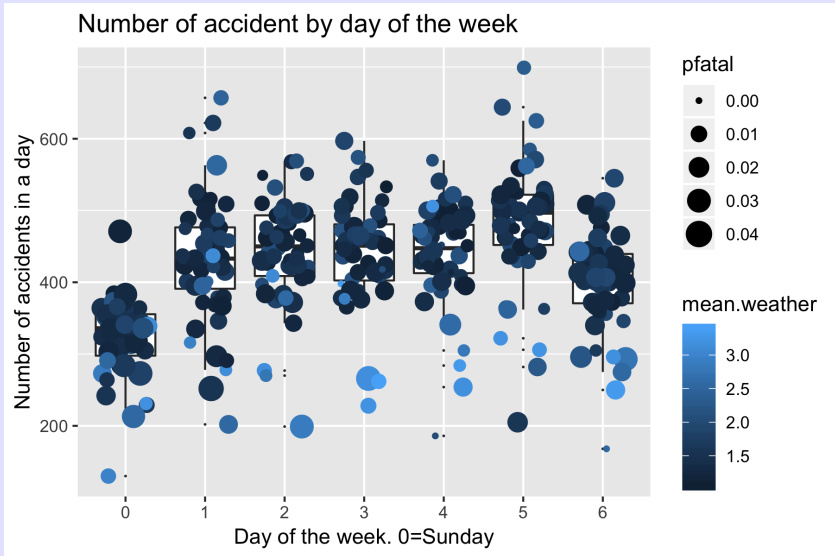
```
1  newplot <- ggplot(data=naccidents,
2                    aes(x=mean.weather, y=freq))+
3    geom_point( ) +
4    geom_smooth(method="loess", color="red", se=FALSE) +
5    ggtitle("Number of accidents by mean weather coded")
6  newplot
```

Number of accidents by mean weather coded

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.

# Split - Apply - Combine - Exercise VII

```
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
```

Number of accident by day of the week

## Split - Apply - Combine - Summary (Simple)

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()*
- The *dplyr* package is specifically design for LARGE data frames and is much faster.

# Split - Apply - Combine - Summary (Simple)

Most simple usage is with *ddply()* and *summarize()*

```
1  sumstat <- plyr::ddply( dataframe, "chunking variable",
2                             plyr::summarize,
3                             v1=....,
4                             v2=....,
5                             v3=...., .... )
```