

# Learning *R*

Carl James Schwarz

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

Reshaping between wide and long formats for data  
Advanced uses.

1. Reshaping between wide and long formats - Advanced
  - 1.1 Wide vs. long formats
  - 1.2 *tidyr* - Melting data frames - wide to long format
  - 1.3 *tidyr* - *spreading* data - long to wide
  - 1.4 *tidyr* - Exercises
  - 1.5 Multiple key-value pairs

## Reshaping Data - Advanced

Wide ↔ Long formats

Wide data format commonly found with many variables or longitudinal data

```
1 > chick.wide <- read.csv("../sampledata/chickweight.csv",
2 +                       header=TRUE, as.is=TRUE,
3 +                       strip.white=TRUE)
4 > head(chick.wide)
```

```
> head(chick.wide)
```

	Chick	Diet	Day01	Day02	Day04	Day06	Day08	Day10	Day12	Day14
1	1	Diet1	42	51	59	64	76	93	106	112
2	2	Diet1	40	49	58	72	84	103	122	131
3	3	Diet1	43	39	55	67	84	99	115	131

We would like a plot of the mean weight over time for each diet.

Long data format transposes each row of data into a long format

```
> head(chick.long)
```

	Chick	Diet	Time	Weight
1	1	1	Day01	42
2	1	1	Day02	51
3	1	1	Day04	59
4	1	1	Day06	64
5	1	1	Day08	76
6	1	1	Day10	93

# Reshaping Data - Why?

- Many statistical models require repeated measure data to be in long format.
- *ggplot()* expects most data to be in long format.
- Quick and dirty way to get plots of multiple variables for screening etc. using faceting in *ggplot()*

- Base *R* *reshape()* function
  - Too hard to use; documentation is useless
- *reshape* - older package do not use
- *reshape2* - deprecated and not updated but still very popular
- *tidyr* - most current but harder to use
- *data.table* - allows for multiple variables to be melted and casted.

```
tidyr::gather(df,  
  key="xxx"  
  value="yyy" c("aaa", "bbb", "ccc", ...))
```

Separate variables into:

- *key* = variable to hold the names of the transposed variables
- *value* = variable to hold the value of the transposed variables
- ... = variable names to transpose (the measure variables)

All other variables (not in the ...) remain in the wide format.



Example:

```
1 chick.long <- tidyr::gather(chick.wide,
2                               key="Time",      # long descriptor
3                               value="Weight",  # long values
4                               c("Day01", "Day02", "Day04",
5                                 "Day06", "Day08", "Day10",
6                                 "Day12", "Day14", "Day16",
7                                 "Day18", "Day20", "Day21"))
8                               )
9 head(chick.long)
```

If you leave out the list of measured variables, then all variables will be transposed.

Example:

```
> head(chick.wide)
```

	Chick	Diet	Day01	Day02	Day04	Day06	Day08	Day10	Day12	Day14
1	1	Diet1	42	51	59	64	76	93	106	120

```
> head(chick.long)
```

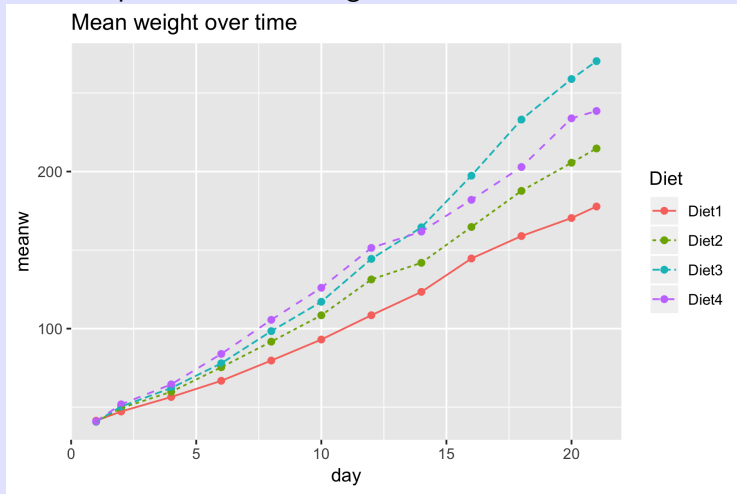
	Chick	Diet	Time	Weight
1	1	1	Day01	42
2	1	1	Day02	51
3	1	1	Day04	59
4	1	1	Day06	64
5	1	1	Day08	76
6	1	1	Day10	93

Now we can compute means for each diet x day combination

```
1 meanw <- plyr::ddply(chick.long, c("Time","Diet"),
2                       summarize,
3                       day = as.numeric(substring(Time[1],4)),
4                       meanw=mean(Weight, na.rm=TRUE))
5 meanw
6
7 plot.meanw <- ggplot2::ggplot(data=meanw,
8                               aes(x=day, y=meanw,
9                                   color=Diet, linetype=Diet))+
10    geom_point()+
11    geom_line()+
12    ggtitle("Mean weight over time")
13 plot.meanw
```

# Plotting with *ggplot2* - Scatterplot

Create a plot of calories vs. grams of fat



Example: melting different variables for plotting purposes:  
Calories from fat, protein, carbohydrates across shelves.

```
1 head(cereal[,c("name","fat","protein","carbo")])
2
3 cereal$calories.fat      <- cereal$fat * 9
4 cereal$calories.protein <- cereal$protein *4
5 cereal$calories.carbo   <- cereal$carbo * 4
6
7 cereal.long <- tidyr::gather(cereal,
8                             key="Source",
9                             value="Calories",
10                            c("calories.fat",
11                              "calories.protein",
12                              "calories.carbo"))
13 head(cereal.long)
```

Example: comparing calories from different sources.

```
> head(cereal[,c("name", "fat", "protein", "carbo")])
```

	name	fat	protein	carbo
1	100%_Bran	1	4	5.0
2	100%_Natural_Bran	5	3	8.0
3				

```
> head(cereal.long)
```

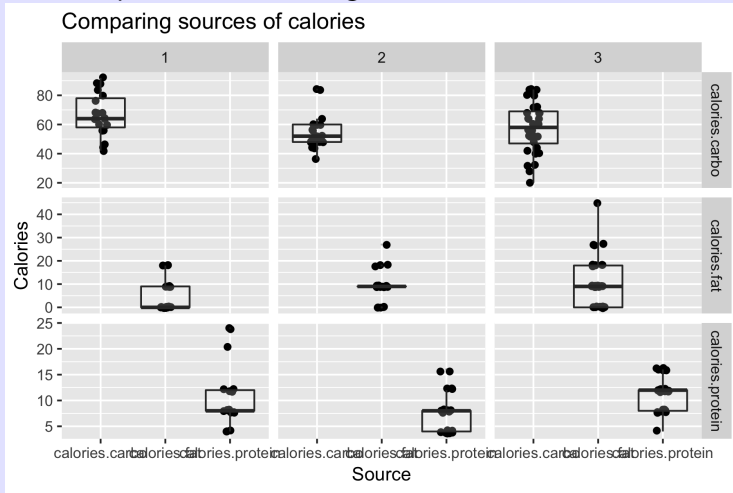
	name	shelfF	Source	Calories
1	100%_Bran	3	calories.fat	9
2	100%_Bran	3	calories.protein	16
3	100%_Bran	3	calories.carbo	20
4	100%_Natural_Bran	3	calories.fat	45
5	100%_Natural_Bran	3	calories.protein	12
6	100%_Natural_Bran	3	calories.carbo	32

Example: comparing calories from different sources.

```
1 plot1 <- ggplot(data=cereal.long, aes(x=Source, y=Calories))
2   ggtitle("Comparing sources of calories")+
3   geom_point(position=position_jitter(w=0.1))+
4   geom_boxplot(alpha=0.2, outlier.size=0)+
5   facet_grid(Source ~ shelfF, scales="free")
6 plot1
```

# Plotting with *ggplot2* - Scatterplot

Create a plot of calories vs. grams of fat





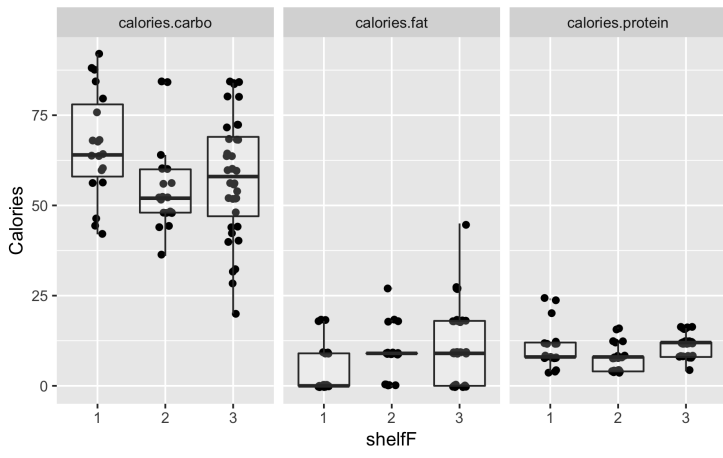
Example: comparing calories from different sources.

```
1 plot2 <- ggplot(data=cereal.long, aes(x=shelfF, y=Calories))
2   ggtitle("Comparing sources of calories")+
3   geom_point(position=position_jitter(w=0.1))+
4   geom_boxplot(alpha=0.2, outlier.size=0)+
5   facet_wrap(~Source )
6 plot2
```

# Plotting with *ggplot2* - Scatterplot

Create a plot of calories vs. grams of fat

Comparing sources of calories



Less common:

```
tidyr::spread(df,  
  key="xxx",  
  value="yyyyy")
```

Example:

```
1 chick.wide2 <- tidyr::spread(chick.long,  
2                               key="Time",  
3                               value="Weight")  
4 head(chick.wide2)  
5
```

Example:

```
> head(chick.long)
```

	Chick	Diet	Time	Weight
1	1	1	Day01	42
2	1	1	Day02	51
3	1	1	Day04	59
4	1	1	Day06	64
5	1	1	Day08	76
6	1	1	Day10	93

```
> head(chick.wide2)
```

	Chick	Diet	Day01	Day02	Day04	Day06	Day08	Day10	Day12	Day14
1	1	Diet1	42	51	59	64	76	93	106	122
2	2	Diet1	40	49	58	72	84	103	122	133

Teeth dataset - number of types of teeth for mammals

- Read the data
- Sum the top and bottom teeth classification.
- Melt the 4 types of teeth
- Make a nice plot comparing the distribution of teeth by mammal classification (H or C)

## Reshaping data - Exercise

```
1 teeth.wide <- read.csv("../sampledata/Teeth.csv",
2                          header=TRUE, strip.white=TRUE,
3                          as.is=TRUE)
4 teeth.wide
5
6 teeth.wide$Incisors <- teeth.wide$Top.incisors + teeth.wide$
7 teeth.wide$Canines <- teeth.wide$Top.canines + teeth.wide$
8 teeth.wide$Premolars <- teeth.wide$Top.premolars + teeth.wide$
9 teeth.wide$Molars <- teeth.wide$Top.molars + teeth.wide$
10
11 head(teeth.wide[,c("Mammal", "Class", "Incisors", "Canines", "Pr

> head(exp.long)
> head(teeth.wide[,c("Mammal", "Class", "Incisors", "Canines", "
      Mammal Class Incisors Canines Premolars Molars
1      BADGER      c         6         2         6         3
2      COUGAR      c         6         2         5         2
3 ELEPHANT SEAL    c         3         2         8         2
```

## Reshaping data - Exercise

```
1 tteeth.long <- tidyr::gather(teeth.wide,  
2                             key="Tooth.Type",  
3                             value="Teeth",  
4                             c("Incisors","Canines","Premolars","Molars"  
5                             )  
6 head(teeth.long)  
7
```

```
> head(teeth.long)
```

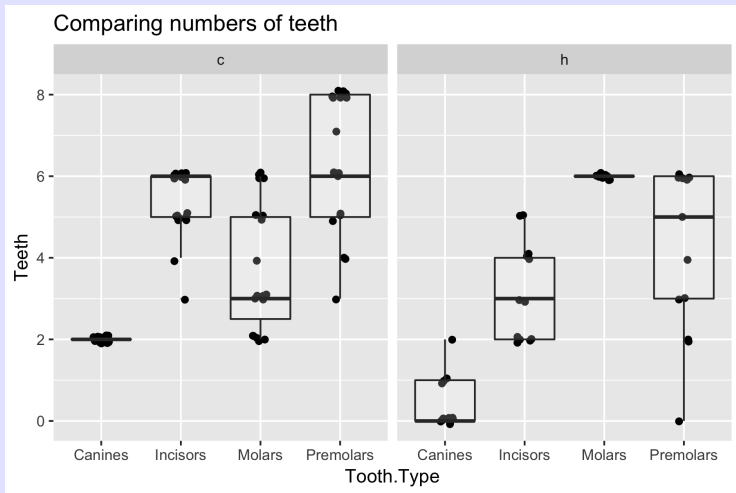
	Mammal	Class	Tooth.Type	Teeth
1	BADGER	c	Incisors	6
2	COUGAR	c	Incisors	6
3	ELEPHANT SEAL	c	Incisors	3

## Reshaping data - Exercise

```
1 plot1 <-ggplot(data=teeth.long, aes(x=Tooth.Type, y=Teeth))-
2   ggtitle("Comparing numbers of teeth")+
3   geom_point(position=position_jitter(h=.1, w=.1))+
4   geom_boxplot(alpha=0.2, outlier.size=0)+
5   facet_wrap(~Class)
6 plot1
7
```



# Reshaping data - Exercise

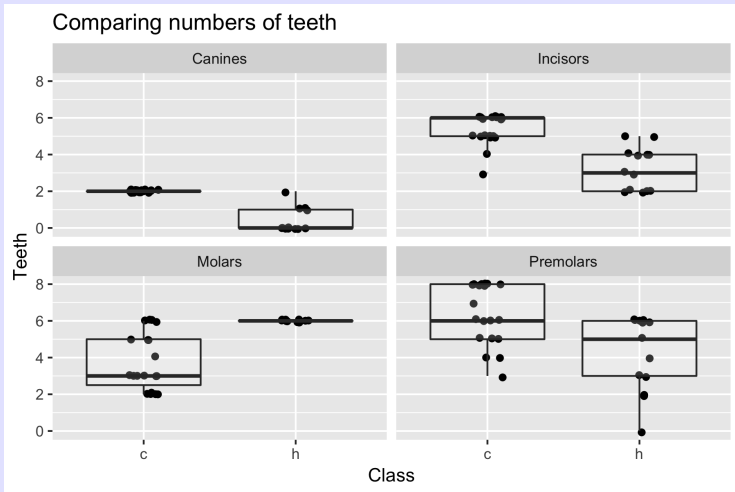


What do you conclude?

## Reshaping data - Exercise

```
1 plot2 <- ggplot(data=teeth.long, aes(x=Class, y=Teeth))+
2   ggtitle("Comparing numbers of teeth")+
3   geom_point(position=position_jitter(h=.1, w=.1))+
4   geom_boxplot(alpha=0.2, outlier.size=0)+
5   facet_wrap(~Tooth.Type, ncol=2)
6 plot2
7
```

# Reshaping data - Exercise



What do you conclude?

Return to the Birds 'n Butts dataset.

Look at the Experimental tab in the Excel file.

Paired design with two halves of 1 nest receiving treatments.

- Read directly in from Excel spreadsheet.
  - How do you skip the first row?
  - How do you fix the variable names?
- Cast to put both values of number of mites on same record
- Compute the difference in the number of mites
- Make a plot to decide if there is an effect?

## Reshaping data - Exercise

```
1 exp.long <- readxl::read_excel("../sampledata/bird-butts-data",
2                               sheet="Experimental", skip=1,
3                               .name_repair="universal")
4 head(exp.long)
5
```

```
> head(exp.long)
```

```
# A tibble: 6 x 5
```

	Nest	Species	Nest.content	Number.of.mites	Treatment
	<dbl>	<chr>	<chr>	<dbl>	<chr>
1	1	HOSP	empty	0	control
2	1	HOSP	empty	0	experimental
3	2	HOSP	empty	1	control
4	2	HOSP	empty	0	experimental

## Reshaping data - Exercise

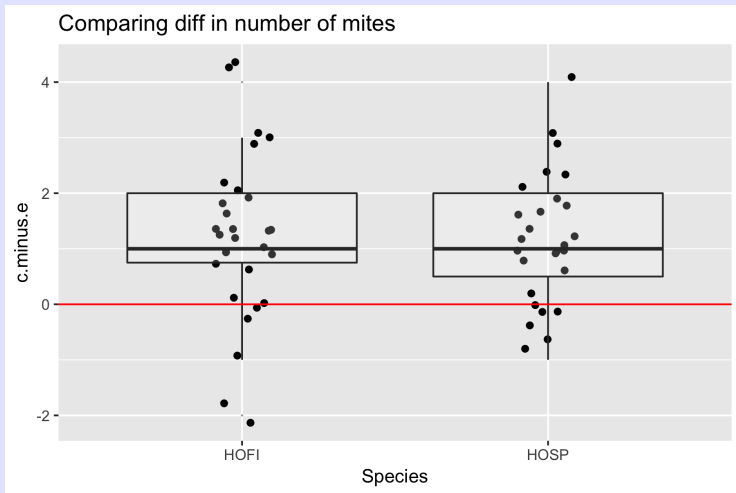
```
1 exp.wide <- tidyr::spread(exp.long,
2                           key="Treatment",
3                           value="Number.of.mites")
4 exp.wide$c.minus.e <- exp.wide$control - exp.wide$experimental
5 head(exp.wide)
6
```

```
> head(exp.wide)
```

	Nest	Species	Nest.content	control	experimental	c.minus.e
1	1	HOSP	empty	0	0	0
2	2	HOSP	empty	1	0	1
3	3	HOSP	eggs	3	1	2

```
1 plot1 <- ggplot(data=exp.wide, aes(x=Species, y=c.minus.e))-
2   ggtitle("Comparing diff in number of mites")+
3   geom_point(position=position_jitter(w=0.1))+
4   geom_boxplot(alpha=0.2, outlier.size=0)+
5   geom_hline(yintercept=0,color="red")
6 plot1
7
```

# Reshaping data - Exercise



What do you conclude?



## Reshaping Data - Advanced

Wide ↔ Long formats

Multiple key-value pair.

*data.table* package.

## Reshaping Data - melting - wide → long

Sometime you have multiple key-value pairs.

Refer to the *ncs\_teaching\_childhealth.xlsx* workbook and the *wide* format worksheet.

```
CHILD_PIDX VISIT_WT_6 VISIT_WT_12 VISIT_WT_18 VISIT_WT_24  
  <chr>      <chr>      <chr>      <chr>      <chr>  
1 a00058528  18          20          21          NA  
2 a00103956  14          NA           24          24  
3 a00104038  11          22          25          NA  
4 a00104358  18          NA           27          31  
5 a00145110  NA          21          24          NA  
6 a00145135  21          NA           31          37  
# É with 10 more variables: CHILD_AGE_18 <chr>, CHILD_AGE_24  
#   GASTRO_12 <chr>, GASTRO_18 <chr>, GASTRO_24 <chr>, EAR_I  
#   EAR_INFECTION_12 <chr>, EAR_INFECTION_18 <chr>, EAR_INF
```

Age, weight, and other variables measured at 4 followup visits and recorded in the wide format.

Example. Very similar to the *reshape2* package.

```
1 library(data.table)
2 child.long <- data.table::melt(as.data.table(child.wide),
3   id.var="CHILD_PIDX",
4   measure.vars=list(c("VISIT_WT_6", "VISIT_WT_12", "VISIT
5                       c("CHILD_AGE_6", "CHILD_AGE
6   variable.name="Visit",
7   value.name=c("Weight", "Age"))
```

Don't forget to use the *as.data.table()* function otherwise you get an uninformative error message.

## Reshaping Data - casting - long → wide

Sometime you have multiple key-value pairs.

Refer to the *ncs\_teaching\_childhealth.xlsx* workbook and the *long* format worksheet.

```
> head(child.long)
# A tibble: 6 x 11
  CHILD_PIDX CHILD_AGE VISIT VISIT_WT CHILD_HEALTH GASTRO D
  <chr>      <dbl> <dbl>   <dbl>   <dbl> <dbl>
1 a00058528     6.7     6     18     1     2
2 a00058528    12.2    12    20     1     1
3 a00058528    17.4    18    21     1    NA
4 a00103956     5.9     6    14     1     2
5 a00103956    18.6    18    24     1    NA
6 a00103956    23     24    24     1    NA
```

Age, weight, and other variables measured at 4 followup visits and recorded in the long format.

Example. Very similar to the *reshape2* package.

```
1 library(data.table)
2 child.wide <- data.table::dcast(as.data.table(child.long),
3   CHILD_PIDX ~ VISIT,
4   value.var=c("VISIT_WT", "CHILD_AGE", "GASTRO", "EAR_INF")
5 head(child.wide)
```

Don't forget to use the *as.data.table()* function, otherwise you get an uninformative error message.

Sometimes necessary to convert from wide ↔ long formats.  
*tidyr* package.

- Wide to long *tidyr::gather()* is often used prior to making a plot using *ggplot()*.
- Long to wide *tidyr::spread()* is often used to bring elements of paired data together.
- *tidyr* has less functionality than *reshape2* but most of the time you don't need the additional features.

I actually prefer the *reshape2* package.

Sometimes necessary to convert from wide ↔ long formats.  
*data.table* package.

- Wide to long *data.table::melt()* is often used prior to making a plot using *ggplot()*.
- Long to wide *data.table::cast()* is often used to bring elements of paired data together.
- *data.table* has more functionality than *reshape2*.
- Don't forget the *as.data.table()* to convert the *data.frame* to a *data.table*.