

Getting back into R - and how to handle messy data

Advanced Research Methods and Skills

2022/03/22

Topics previously covered

- Data visualization
 - Using **ggplot2**
- Data manipulation
 - Using **dplyr**
- Basic statistics
 - t-tests, correlations
 - regression
 - ANOVA

Themes for today

- Refamiliarizing yourself with R
- Basic data transformations
- Missing data

The screenshot shows the Studio Cloud interface. The 'New Project' button is circled in red. Below the screenshot is a diagram illustrating data structure concepts using a table of COVID-19 data.

country	year	cases	population
Afghanistan	2020	15	170071
Afghanistan	2020	296	205360
Brazil	2020	3157	17256362
Brazil	2020	8168	17434898
China	2020	21398	12725272
China	2020	1286	1280583

The diagram illustrates three ways to view the data:

- variables:** Vertical arrows point to the columns of the table, representing individual variables.
- observations:** Horizontal arrows point to the rows of the table, representing individual observations.
- values:** Circles are placed at the intersection of each row and column, representing individual data values.

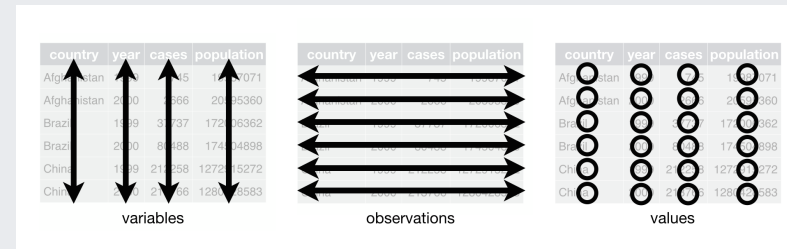
Tabular data

```
## # A tibble: 16 x 4
##   Participant Viewpoint Block    RT
##   <int> <chr> <chr> <dbl>
## 1         1    1 Different First    516.
## 2         2    1 Different Second  409.
## 3         3    1 Same      First    419.
## 4         4    1 Same      Second  414.
## 5         5    2 Different First    494.
## 6         6    2 Different Second  385.
## 7         7    2 Same      First    516.
## 8         8    2 Same      Second  385.
## 9         9    3 Different First    507.
## 10        10    3 Different Second  371.
## 11        11    3 Same      First    475.
## 12        12    3 Same      Second  374.
## 13        13    4 Different First    470.
## 14        14    4 Different Second  394.
## 15        15    4 Same      First    434.
## 16        16    4 Same      Second  414.
```

Tables of data are what you're most commonly dealing with in R.

This one conforms to the **tidy data** principles:

One row per observation, one column per variable



Different types of file

The most common file formats you'll deal with are either Excel files or text files, but you may also find dealing with SPSS files useful.

Fortunately, R has several functions and packages for importing data!

File formats	File extension	Functions	Package
SPSS	.sav	<code>read_sav()</code>	<code>library(haven)</code>
Excel	.xls, .xlsx	<code>read_excel()</code>	<code>library(readxl)</code>
Text	.csv, .txt, .*	<code>read_csv()</code> , <code>read_delim()</code>	<code>library(readr)</code>

Data wrangling

`dp`lyr is a really useful package for manipulation of data tables.

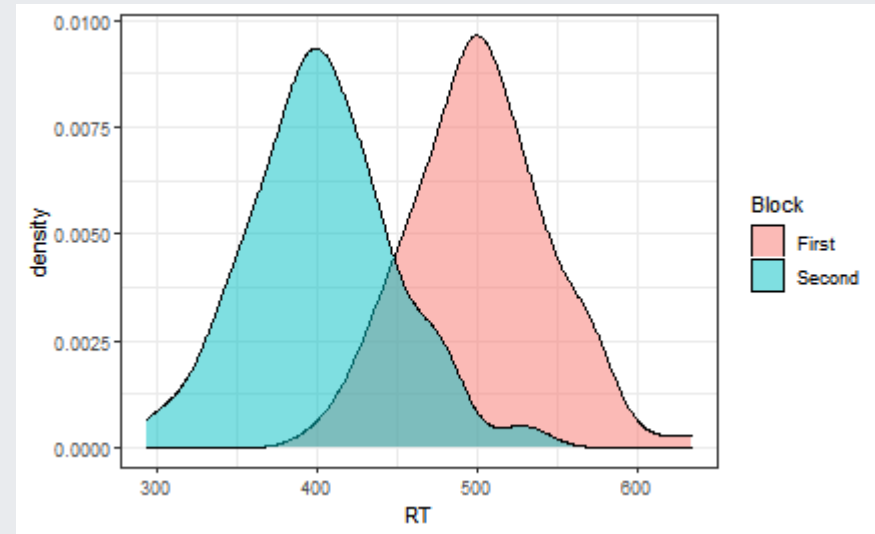
Function	Effect
<code>select()</code>	Include or exclude variables (columns)
<code>arrange()</code>	Change the order of observations (rows)
<code>filter()</code>	Include or exclude observations (rows)
<code>mutate()</code>	Create new variables (columns)
<code>group_by()</code>	Create groups of observations
<code>summarise()</code>	Aggregate or summarise groups of observations (rows)

Quickly plotting your data

Plotting data is really important for all aspects of data analysis.

The `ggplot2` package provides a framework for doing this:

```
ggplot(example_rt_df,  
       aes(x = RT, fill = Block)) +  
  geom_density(alpha = 0.5) +  
  theme_bw()
```



Running statistics

The humble t-test...

```
t.test(RT~Block, data = example_rt_df)
```

```
##  
##      Welch Two Sample t-test  
##  
## data:  RT by Block  
## t = 15.939, df = 197.45, p-value < 2.2e-16  
## alternative hypothesis: true difference in means between group First and group Second is not equal  
## 95 percent confidence interval:  
##    86.71124 111.19806  
## sample estimates:  
##  mean in group First mean in group Second  
##           501.7161           402.7615
```


Running statistics

The factorial ANOVA... (the `aov_ez()` function from the `afex` package)

```
aov_ez(dv = "RT",  
       within = c("Block", "Viewpoint"),  
       id = "Participant",  
       data = example_rt_df)
```

```
## Anova Table (Type 3 tests)
```

```
##
```

```
## Response: RT
```

```
##           Effect      df      MSE          F ges p.value  
## 1           Block 1, 49 1627.81 300.77 *** .565 <.001  
## 2           Viewpoint 1, 49 1241.33      3.50 + .011 .067  
## 3 Block:Viewpoint 1, 49 2513.64      0.21 .001 .651
```

```
## ---
```

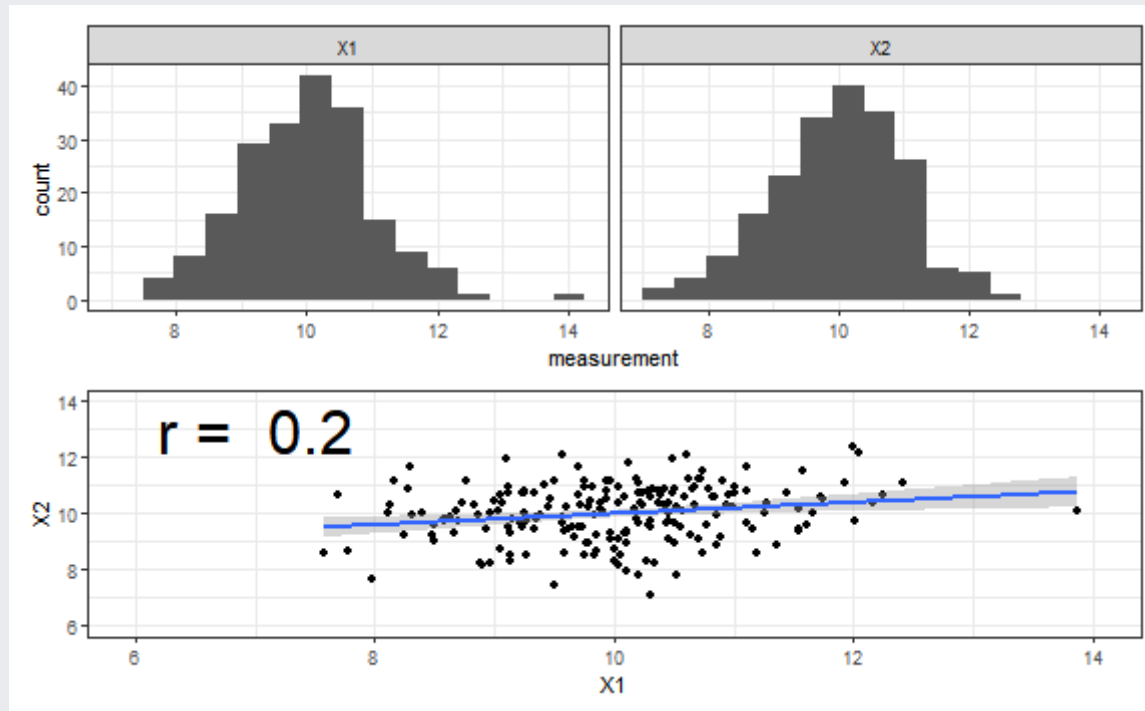
```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '+' 0.1 ' ' 1
```

How to handle "messy" or otherwise awkward data

The ideal data

In an ideal world all our data would be beautifully normal:

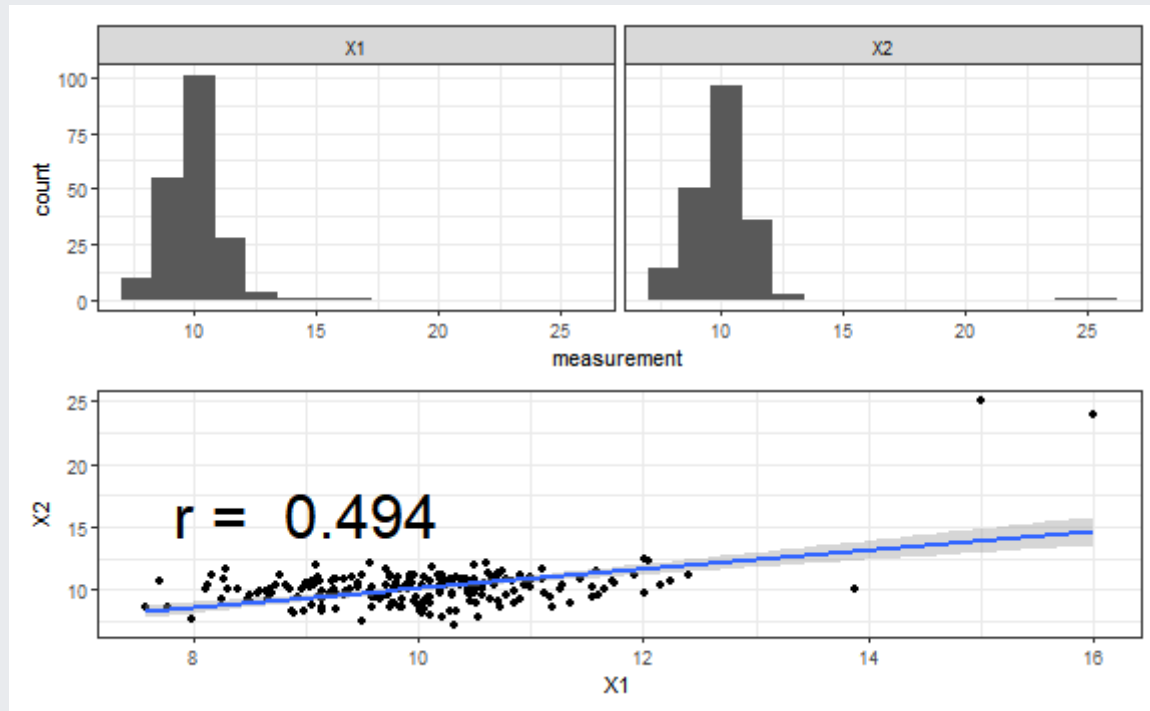
```
## `geom_smooth()` using formula 'y ~ x'
```



The real data

But reality is rarely so kind. Data can be all kinds of messy. It can have *outliers*...

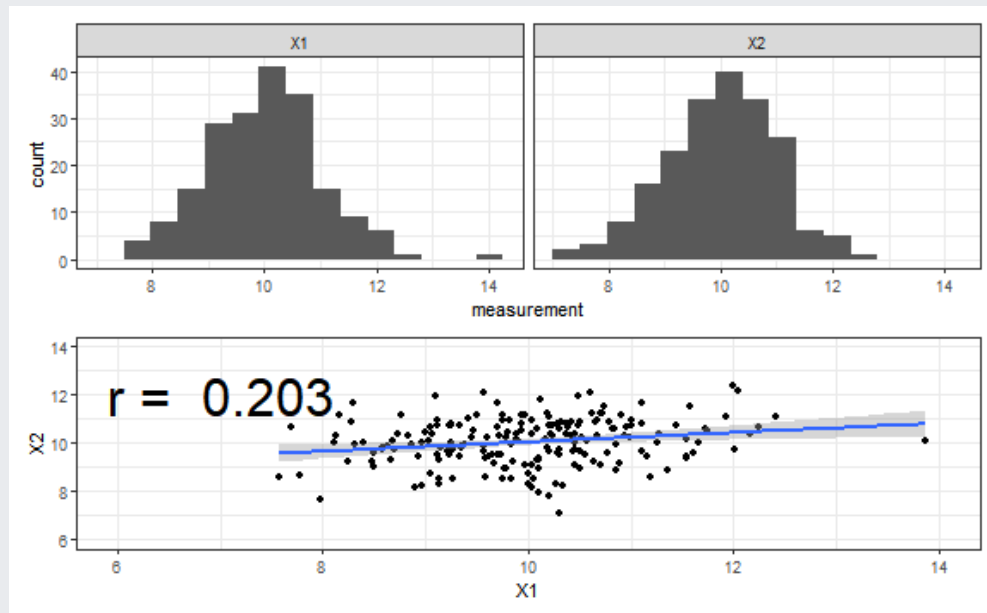
```
## `geom_smooth()` using formula 'y ~ x'
```



The real data

Data can be *missing*...

```
## `geom_smooth()` using formula 'y ~ x'
```

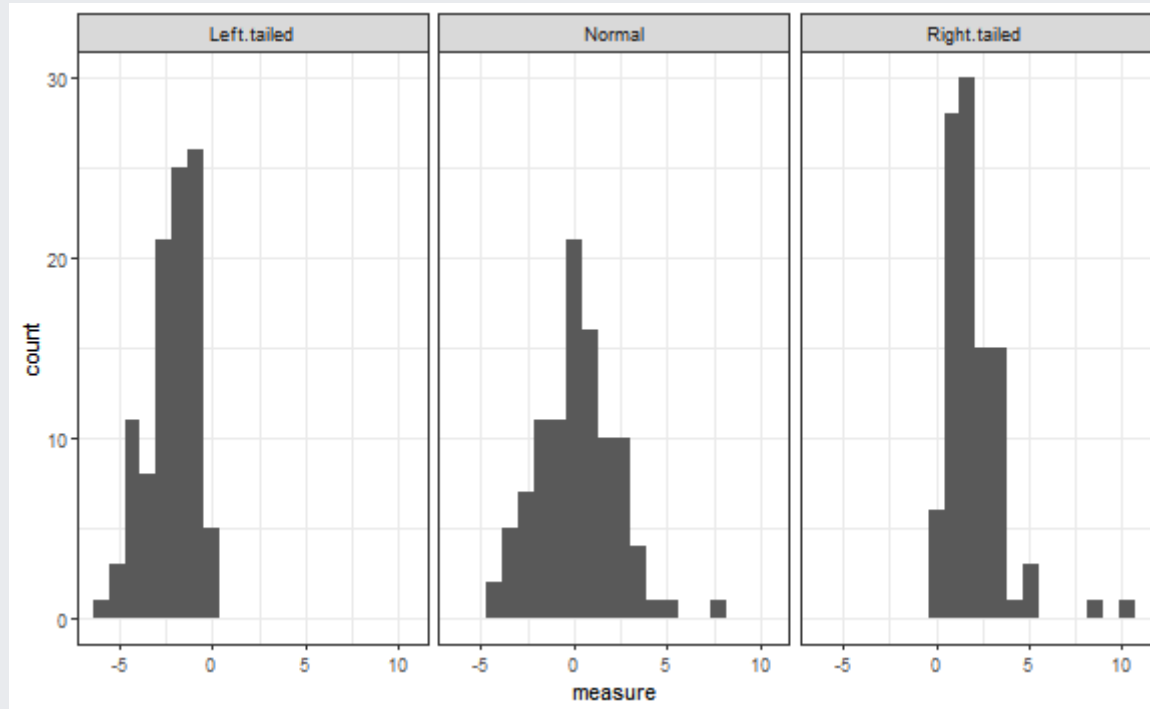


```
##           X1           X2
## 1           NA      8.589878
## 2           NA      8.666062
## 3           NA      7.476062
## 4           NA      8.254748
## 5           NA      9.341976
## 6  10.528092           NA
## 7  10.586026           NA
## 8  10.300834      9.752974
## 9  10.035375     10.365672
## 10  8.896707      8.153483
```

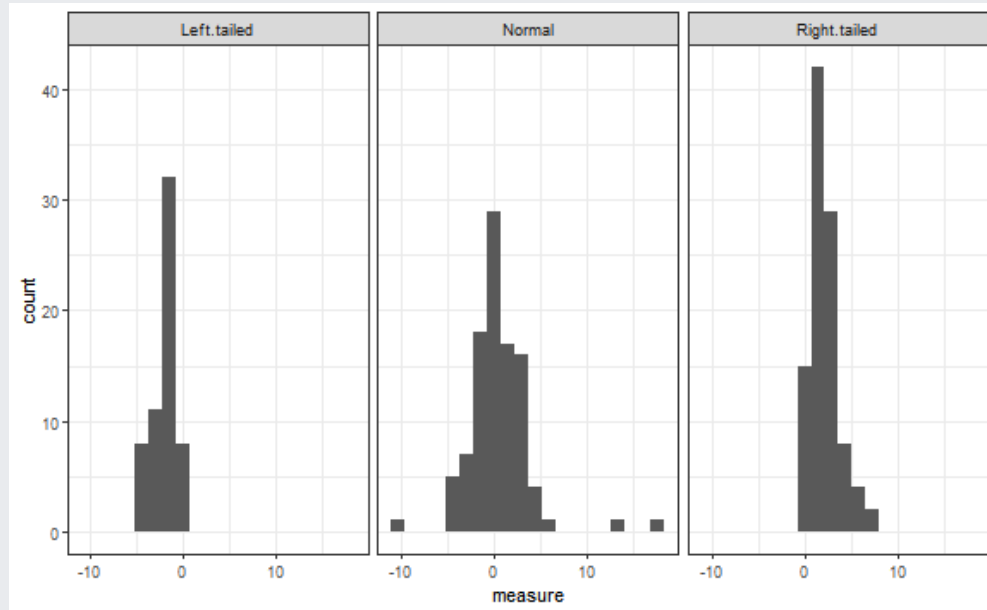
Complete cases = 193

The real data

Data can be *skewed*...



The real data



There can be any combination of these things...

All of these pose problems for estimating the properties of our data, the relationships between variables, and the phenomena we are investigating.

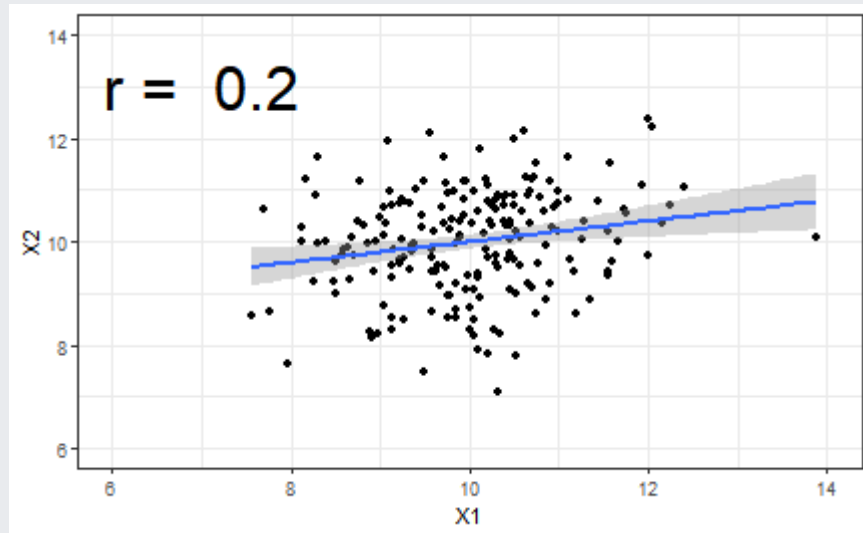
Handling outliers

What is an outlier?

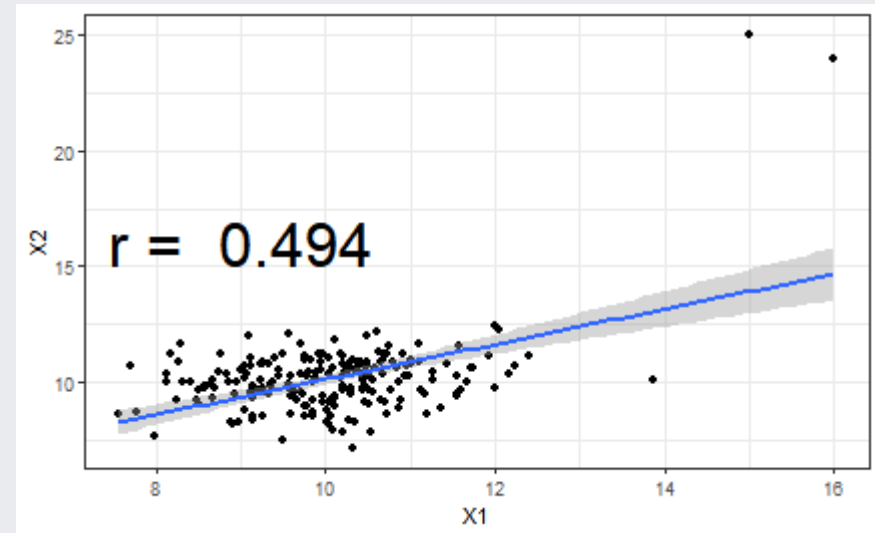
Two out of the 200 pairs of x-y values were replaced.

The resulting coefficient (approx $r = .49$) is *way-off* the true coefficient for these data.

```
## `geom_smooth()` using formula 'y ~ x'
```



```
## `geom_smooth()` using formula 'y ~ x'
```



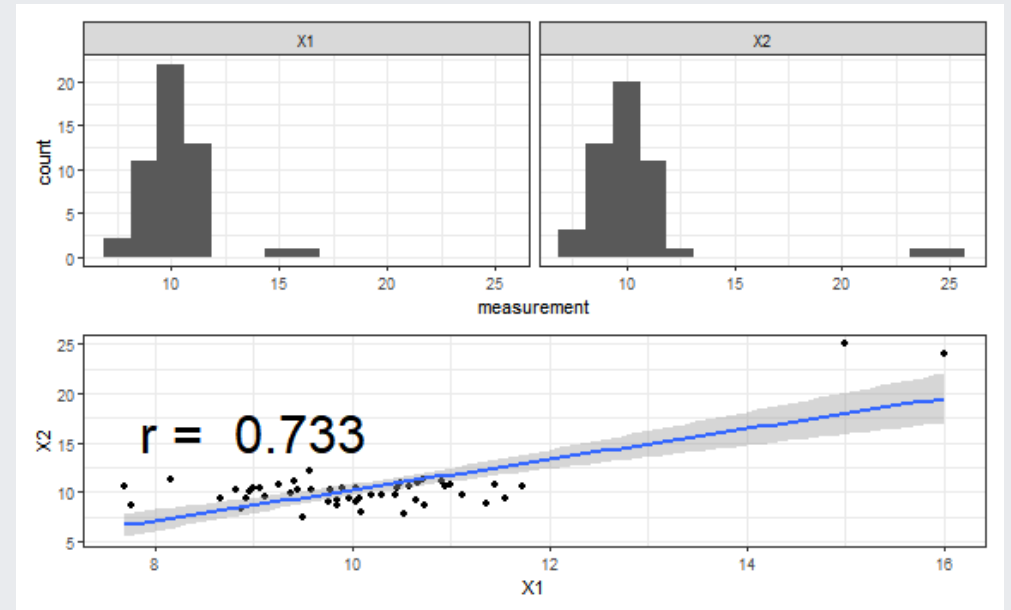
What is an outlier?

The problem gets even worse with smaller sample sizes.

Here there are 50 datapoints with two outliers, rather than 200.

The correlation coefficient becomes even *more* biased than it was previously.

```
## `geom_smooth()` using formula 'y ~ x'
```



What should we do with outliers?

Three common approaches:

1. Remove them

- If you're sure these reflect an error, not genuine data, then removal is a possibility.

2. Transformation

- *rescaling* or *transforming* your data may help reduce the influence of outliers. (We'll come back to this!)

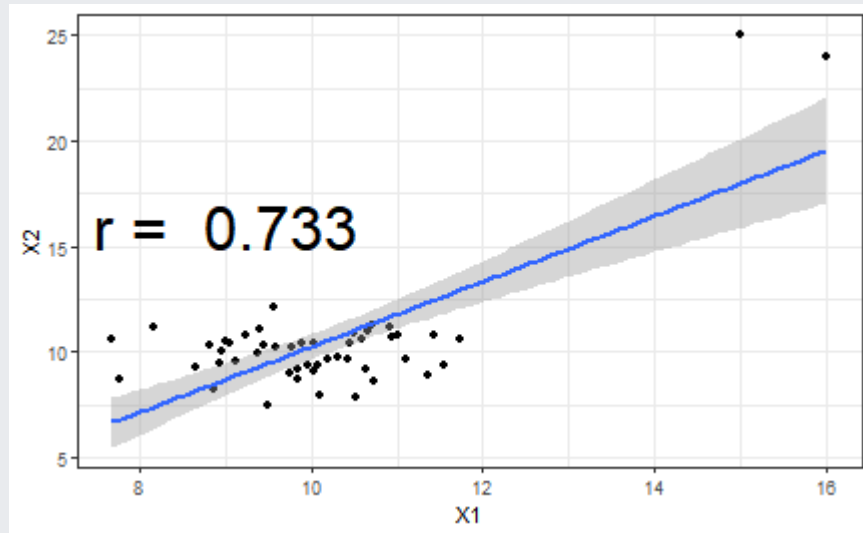
3. Replace them

- Replacing the outliers with values ± 2 -3 standard deviations away from the mean.

Identifying and replacing outliers

Identifying outliers

```
## `geom_smooth()` using formula 'y ~ x'
```

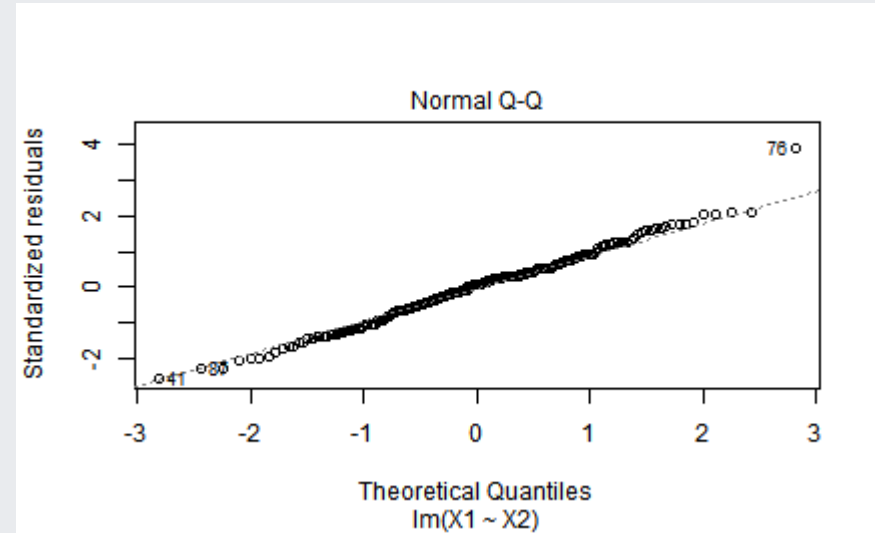
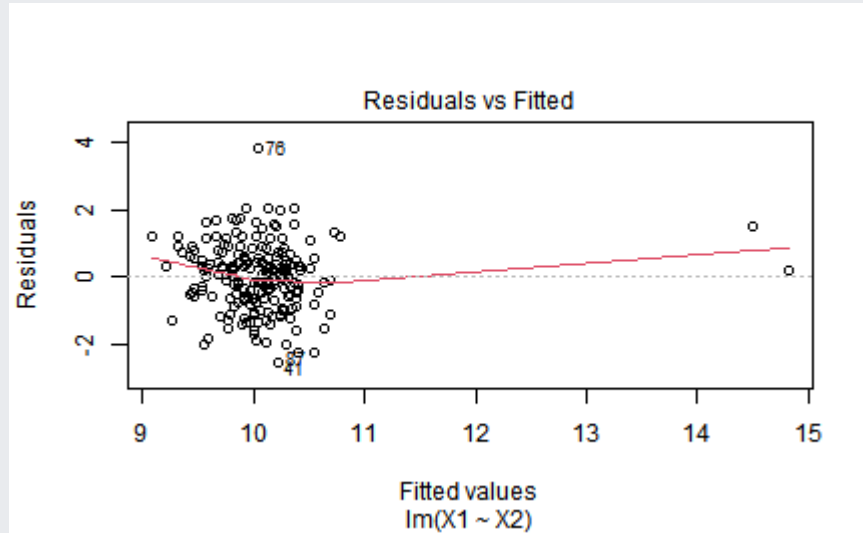


Plotting your data can be an excellent way to spot outliers: here they're *very* obvious!

Identifying outliers

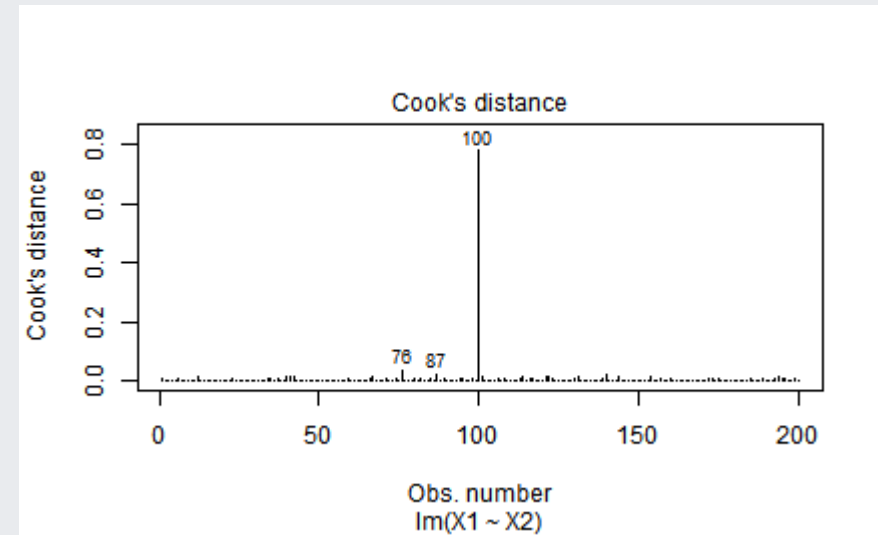
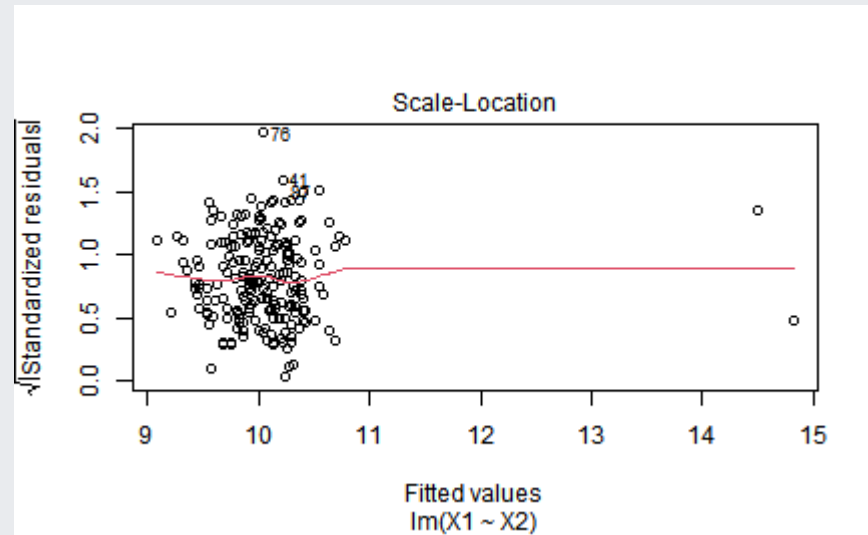
Plotting the residuals of your linear model will also help you identify troublesome observations.

```
plot(lm(X1 ~ X2, data = temp_df_out))
```



Identifying outliers

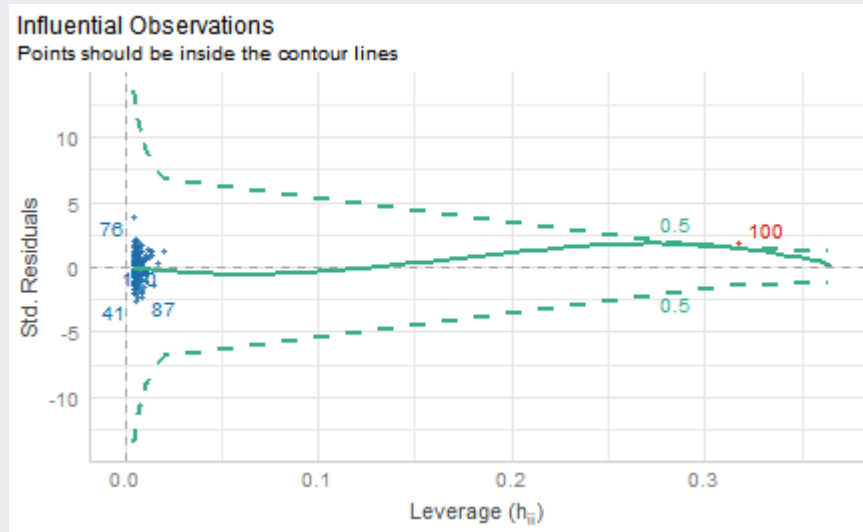
Plotting the residuals of your linear model will also help you identify troublesome observations.



Identifying outliers

The performance library has a `check_outliers()` function that helps too!

```
library(performance)
plot(check_outliers(lm(X1 ~ X2, data = temp_df_out)))
```

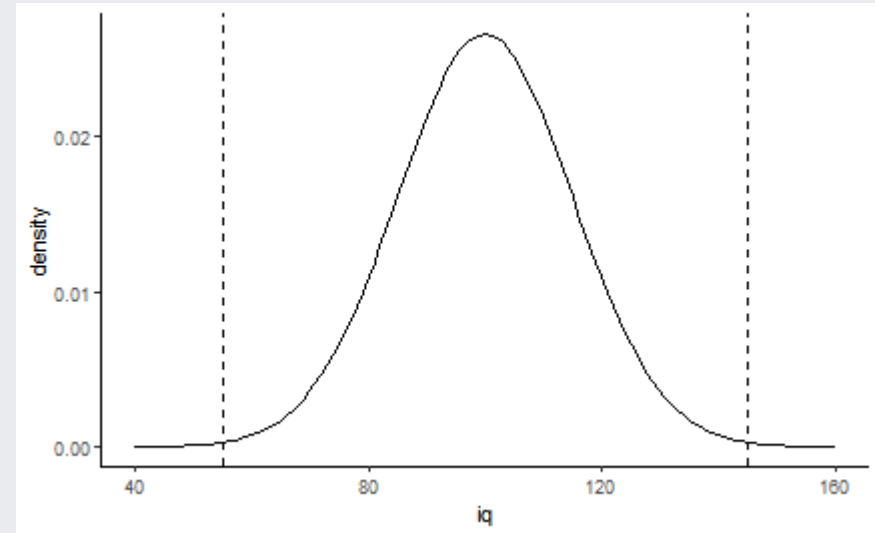


Identifying outliers

Sometimes a *threshold* is used to determine whether an observation is an outlier.

Sometimes this is driven by common sense: e.g. a value of 120 for a participant's age is **extremely** unlikely to be genuine.

Sometimes this is *data-driven*: e.g. values more than ± 3 standard deviations away from the mean are *unusual*.



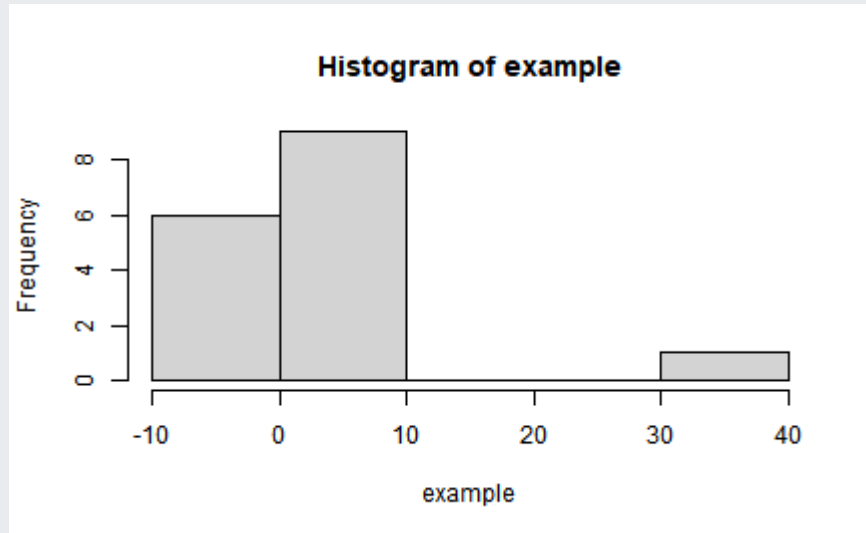
Scaling and standardizing

The data

Manually scale

scale()

```
example <- c(rnorm(15), 35)  
hist(example)
```



Scaling and standardizing

The data

Manually scale

scale()

```
example
```

```
## [1] 0.54240596 1.73336700 0.52365537 1.16488774 0.52106485 0.44979269 0.11852738 -0.72022
## [9] -0.14502293 -1.91319909 -1.80588221 2.08138700 -0.05824803 -0.97010176 0.08841495 35.00000
```

```
(example - mean(example)) / sd(example)
```

```
## [1] -0.19857188 -0.06310656 -0.20070465 -0.12776797 -0.20099931 -0.20910613 -0.24678575 -0.34218
## [9] -0.27676316 -0.47788355 -0.46567684 -0.02352119 -0.26689299 -0.37061137 -0.25021088 3.72079
```

Scaling and standardizing

The data

Manually scale

scale()

```
# you don't need to use t() - that transposes the values so they fit on the slide :)  
# Just use scale()  
t(scale(example))
```

```
##           [,1]           [,2]           [,3]           [,4]           [,5]           [,6]           [,7]           [,8]  
## [1,] -0.1985719 -0.06310656 -0.2007047 -0.127768 -0.2009993 -0.2091061 -0.2467858 -0.3421891  
##           [,9]          [,10]          [,11]          [,12]          [,13]          [,14]          [,15]          [,16]  
## [1,] -0.2767632 -0.4778835 -0.4656768 -0.02352119 -0.266893 -0.3706114 -0.2502109 3.720791  
## attr("scaled:center")  
## [1] 2.288177  
## attr("scaled:scale")  
## [1] 8.791631
```

Removing values above a threshold

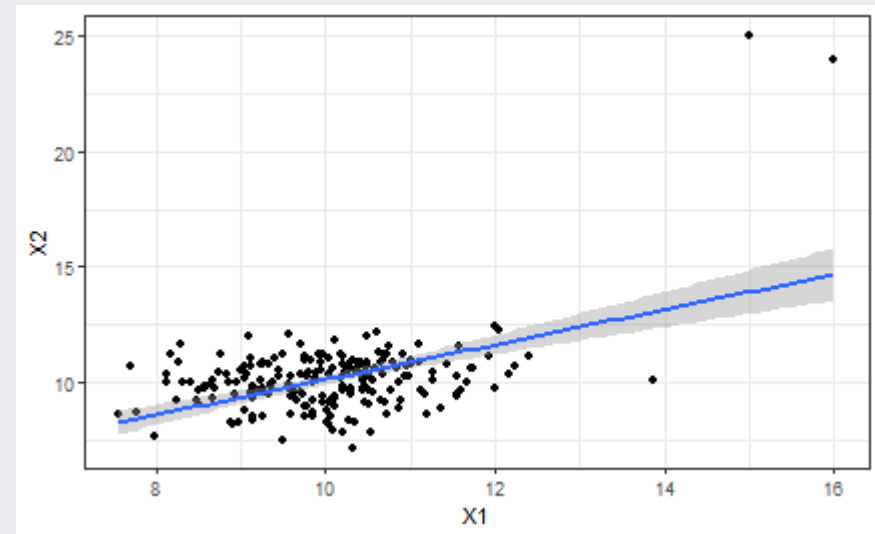
The `filter()` function from `dplyr` can be used to remove outliers easily!

With outliers

Without outliers

```
temp_df_out %>%  
  ggplot(aes(x = X1, y = X2)) +  
  geom_point() +  
  theme_bw() +  
  stat_smooth(method = "lm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```



Removing values above a threshold

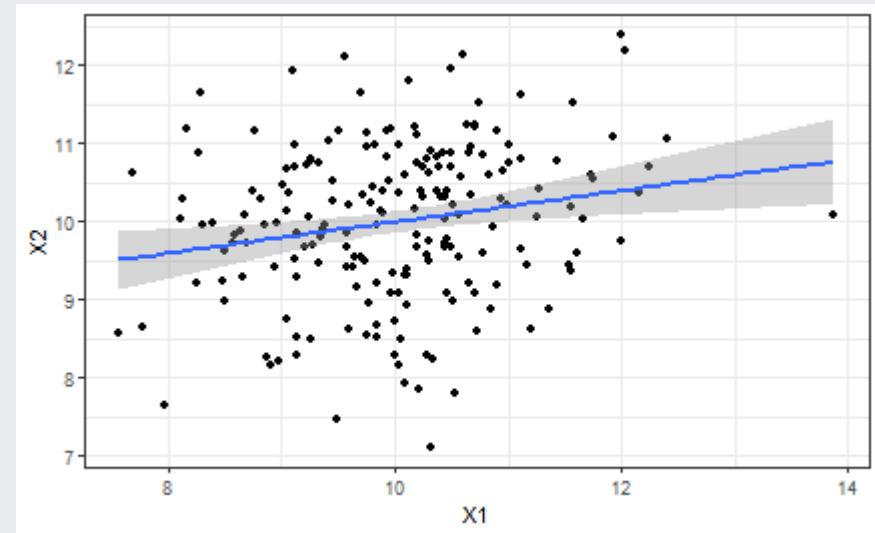
The `filter()` function from `dplyr` can be used to remove outliers easily!

With outliers

Without outliers

```
temp_df_out %>%  
  dplyr::filter(X1 < 15) %>%  
  ggplot(aes(x = X1, y = X2)) +  
  geom_point() + theme_bw() +  
  stat_smooth(method = "lm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```

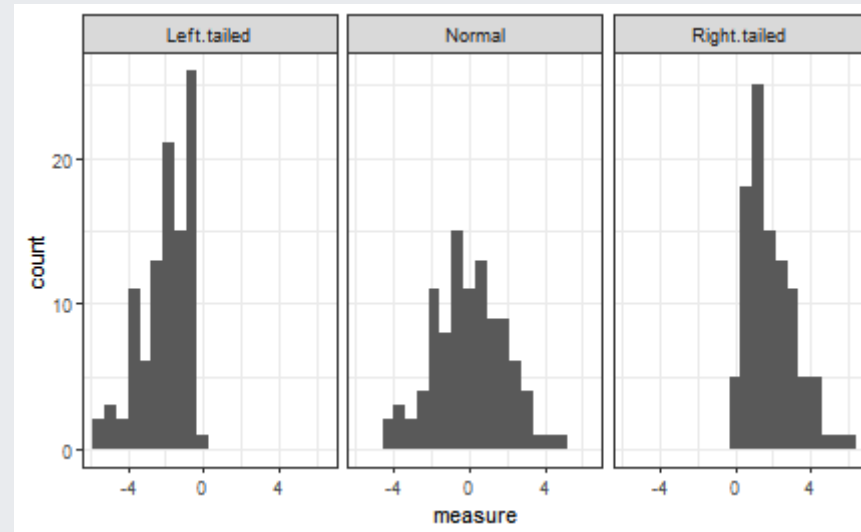


Data transformation

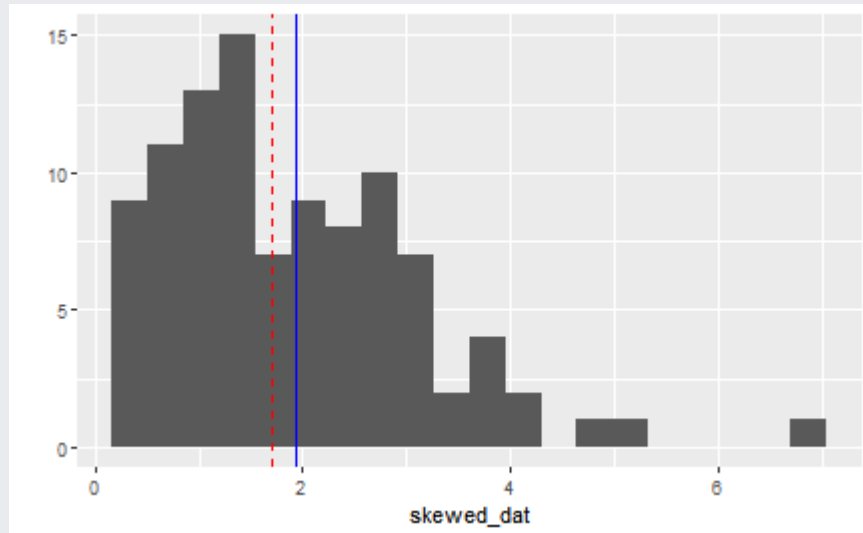
Skewed data

Skewed data is data that *leans* in a particular direction.

These are often described by the direction of the "long-tail" - so a left-tailed distribution means a distribution with a long tail on the left, rather than most values on the left.



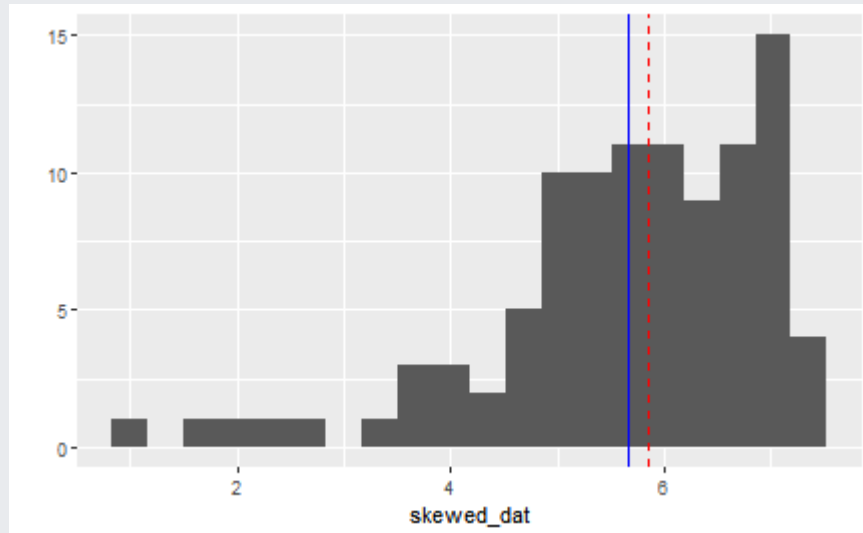
Skewed data



This data is *right-tailed*. This is sometimes also called *positively skewed*. For this type of data, the mean (blue line) is usually higher than the median (red, dashed line).

This type of skew is relatively common with data that is *bounded* at zero. e.g. reaction time data, the distribution of wages

Left-tailed skew



Data skewed the opposite way - many high scores but few low scores - has a long *left* tail. This is also called *negative skew*. The mean (blue solid line) is usually less than the median (red dashed line).

Transformation of skewed data

One way to handle skew is to transform the data to a different *scale*.

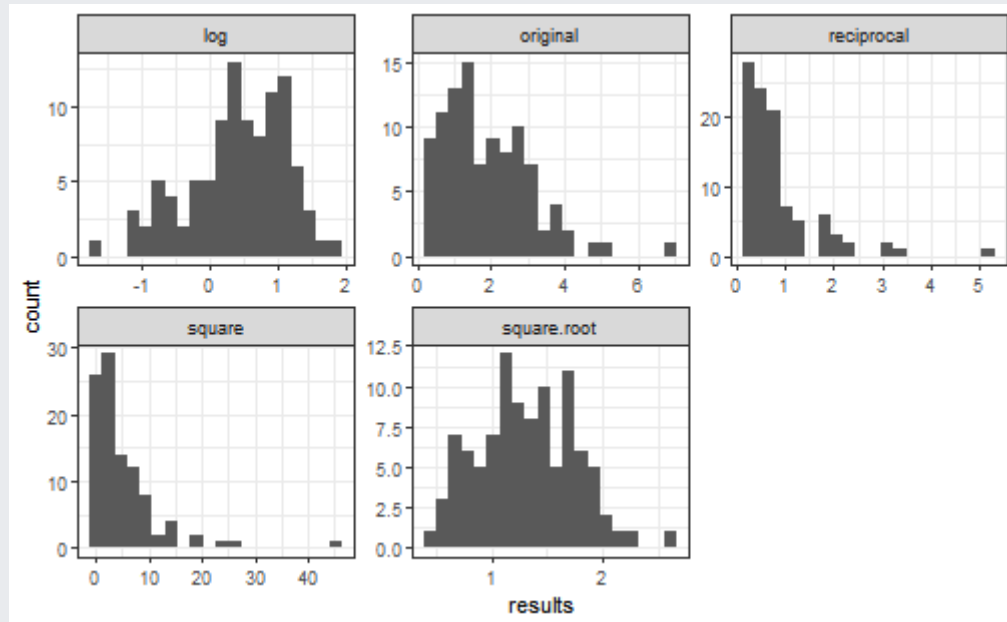
Transformation type	code
Log	$\log(X)$
Square root	$\text{sqrt}(X)$
Reciprocal	$1 / X$
Square	x^2

(See Section 5.8.2 in Field et al., DSUR)

Transformation of skewed data

Right-tailed

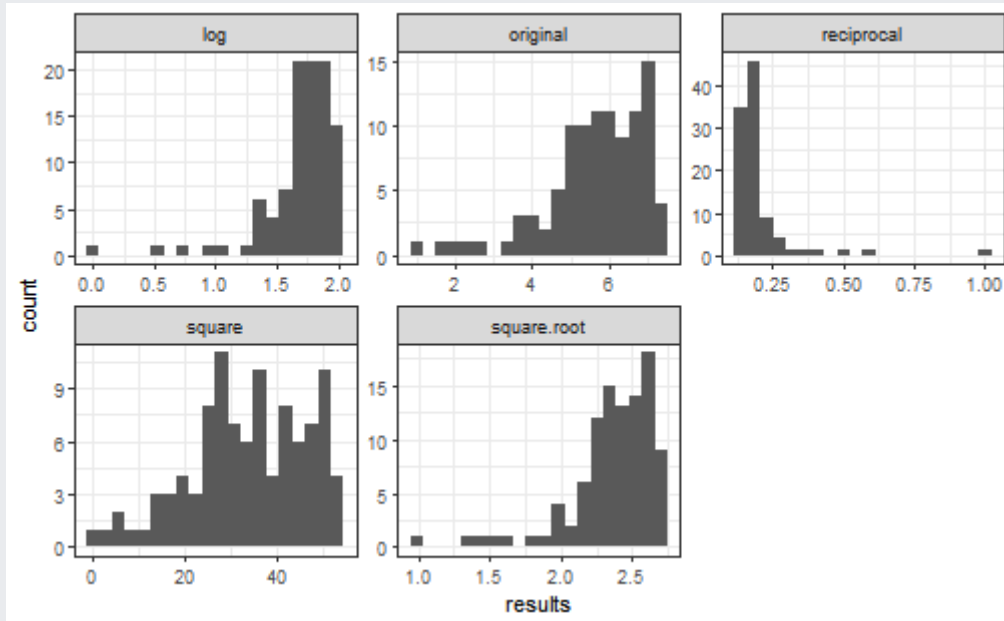
Left-tailed



Transformation of skewed data

Right-tailed

Left-tailed



Handling missing data

Types of missing data

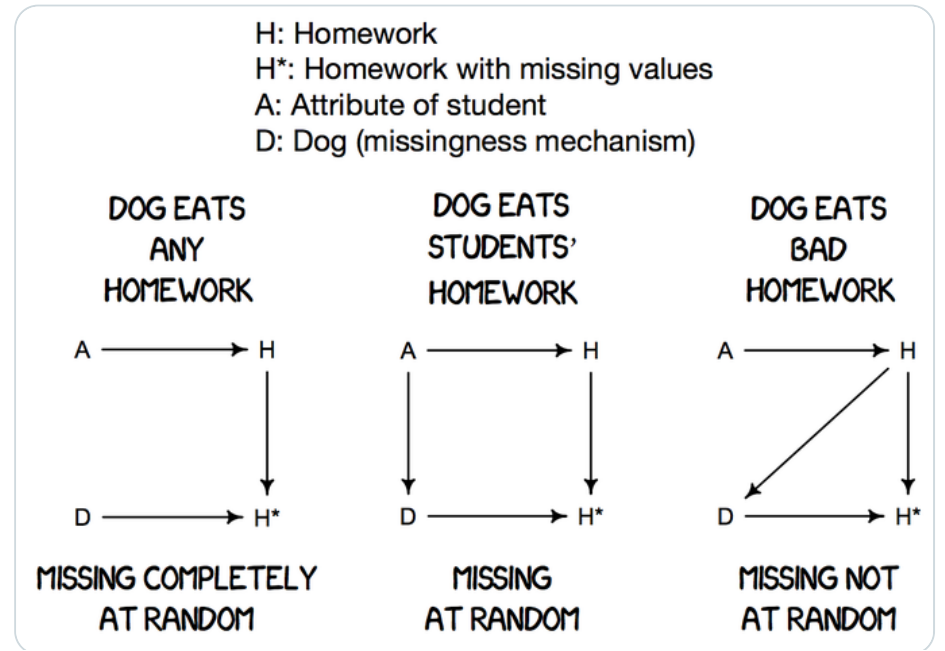
- Missing Completely At Random
 - Missingness does not depend on anything
- Missing At Random
 - Missingness depends on the observed data
- Missing Not At Random
 - Missingness depends on the missing data



Richard McElreath 🐶
@rilmcelreath



In today's lecture, I tried to redefine missing data types (MCAR, MAR, MNAR) as different reasons a dog might eat your homework. This needs more work, but audience seemed to appreciate it.



Missing Completely At Random

If you have missing data, MCAR is the best kind of missing data.

There is nothing *systematic* about which data is missing.

For example, all your participants filled out three different questionnaires. Unfortunately, your dog chewed through a pile of them, and half of your participants now have only two questionnaires.

##		X1	X2	X3
## 1		NA	10.564509	12.00366
## 2		NA	11.134840	12.46372
## 3		NA	8.197005	11.06698
## 4		NA	9.704088	14.19018
## 5		NA	8.706921	13.03839
## 6	8.344476	9.537686	12.89438	
## 7	11.263166	7.791435	13.76203	
## 8	8.669327	8.276546	12.92337	
## 9	8.377079	9.164257	11.07893	
## 10	9.288548	6.947327	12.80245	

Missing At Random

##		X1	X2	X3	age
## 1		11.047154	9.130396	12.86073	19
## 2		8.509500	10.760188	12.07602	20
## 3		11.619125	9.636584	13.04400	19
## 4		11.212559	7.600109	11.90494	20
## 5		10.951217	9.341155	13.35871	19
## 6		9.286955	10.227833	12.48842	19
## 7		9.555602	8.320163	NA	36
## 8		10.394811	7.726169	15.13988	19
## 9		8.376943	8.289583	NA	26
## 10		10.528843	8.812628	13.50900	20

Confusingly, Missing At Random (MAR) data is not missing (completely) at random.

For example, for some reason, people older than 21 typically failed to complete the third questionnaire.

This data is MAR - whether the data in the third column is missing is related to the value of the fourth column.

Missing Not At Random

##		X1	X2	X3	age
## 1		11.095201	8.235348	NA	19
## 2		11.159246	9.198542	NA	19
## 3		11.026114	11.395065	NA	29
## 4		10.723638	10.066800	NA	19
## 5		11.399608	11.517779	NA	19
## 6		10.051001	9.103957	12.61627	18
## 7		9.073119	7.654592	13.10197	42
## 8		10.110520	10.048114	12.38895	29
## 9		9.729318	10.017504	NA	19
## 10		11.901271	9.132584	13.14818	52

The final, most troubling type of missing data is data that is Missing Not At Random (MNAR).

For example, imagine that the questionnaire relates to depression; people who score high for depression are less likely to complete the final questionnaire.

In this case, the values that are missing for the third questionnaire depends on the value of the responses to that questionnaire, so this data is MNAR.

Dealing with missing data

List-wise deletion: Cases with missing data are completely **removed** from **all** analysis.

Pair-wise deletion: Cases with missing data are only **removed** from **comparisons where one or more variables are missing**.

By default, functions such as `mean()` return NA if any value in the input is NA/missing.

```
mean(temp_df_missing$X1)
```

```
## [1] NA
```

```
mean(temp_df_missing$X1, na.rm = TRUE)
```

```
## [1] 10.00553
```

```
sum(complete.cases(temp_df_missing))
```

```
## [1] 193
```

Single Imputation

Replace missing values with a simple "best-guess". e.g. Using the mean or the median for the condition.

```
orig_data <- 1:12  
orig_data
```

```
## [1] 1 2 3 4 5 6 7 8 9 10 11 12
```

```
mean(orig_data)
```

```
## [1] 6.5
```

```
missing_one <- orig_data  
missing_one[6] <- NA  
missing_one
```

```
## [1] 1 2 3 4 5 NA 7 8 9 10 11 12
```

```
mean(missing_one, na.rm = TRUE)
```

```
## [1] 6.545455
```

Single Imputation

Replace missing values with a simple "best-guess". e.g. Using the mean or the median for the condition.

```
replace_one <- missing_one
replace_one[6] <- mean(missing_one,
                      na.rm = TRUE)
mean(replace_one)
```

```
## [1] 6.545455
```

```
mean(orig_data)
```

```
## [1] 6.5
```

Problem: the mean and median are biased by the missing data. And replacing a missing value with one of these values tends to artificially reduce variability.

Multiple Imputation

In **multiple** imputation, we replace missing values with estimates based on a *model* of the data that incorporates uncertainty about what the value should be.

We create a model based on the data that is not missing, and use its predictions to guess the values that the missing data has.

We do this multiple times and then take an average or *pool* the results to fill in the gap.

Packages such as `mice` and `Amelia` can do this for us, and help us identify patterns of missingness.

Alternative approaches to missing data, skew, and other oddities

Generalized Linear Models (as opposed to General Linear Models) allow modelling of data of many different types without necessitating transformations.

For example, counts can be modelled using Poisson regression, and categorical outcomes can be modelled with logistic regression.

Multilevel or *mixed*-models can handle all of these things and much more besides; they are perfectly capable of handling missing data.

We'll cover both logistic regression and multilevel models later in the course!

Next week

Look into power and effect sizes:

See Field et al, Discovering Statistics Using R, pages 56-59, Sections on:

- Type I and Type II error (2.6.3)
- effect sizes (2.6.4)
- statistical power (2.6.5)

Cohen, J. (1992). A power primer. *Psychological Bulletin*, 112(1), 155-159. <http://dx.doi.org/10.1037/0033-2909.112.1.155>