

3 - Data Exploration

Felipe Melo

Nottingham Trent University - UK

You should know today

- Make questions to your data?
- Explore the basic features of your data
- Make simple exploratory graphics

Before we begin



- R and Rstudio installed
- Don't panic
- Everything is reproducible
- You'll have to train to fix the content

What questions should I make to the data?

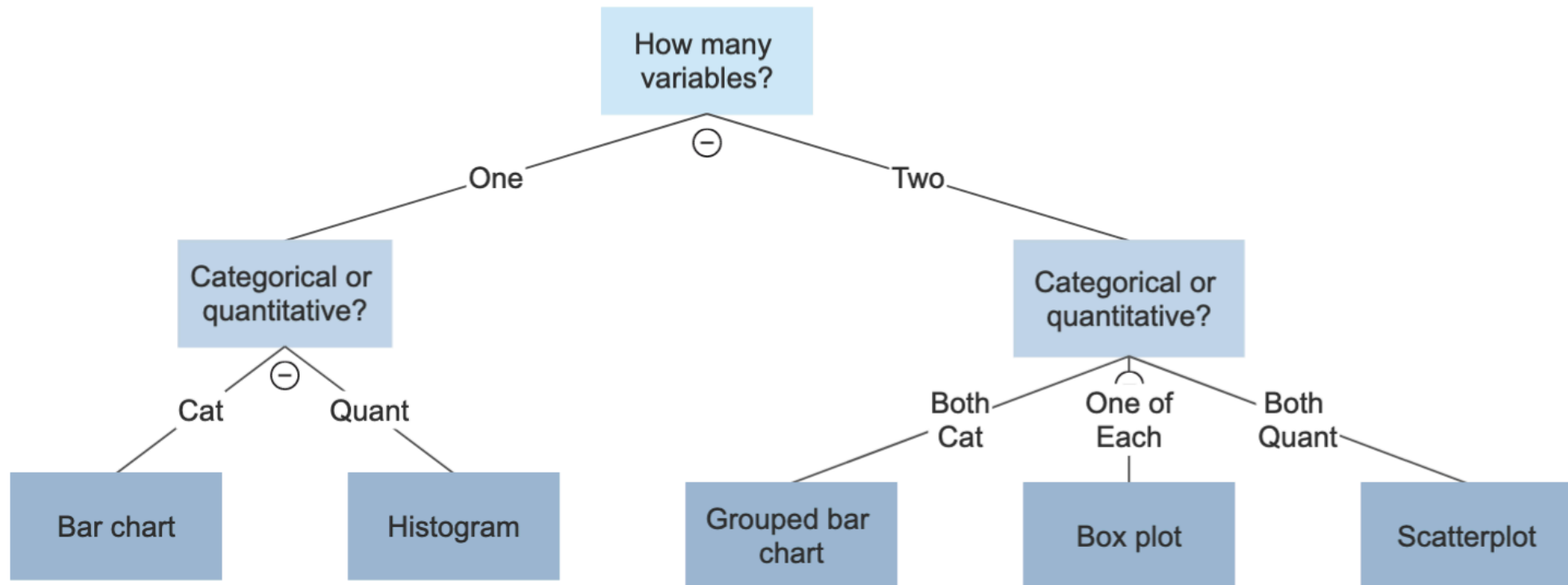
Back to Spreadsheets

The Penguins file

```
1 penguins_df<-read.csv("https://raw.githubusercontent.com/fplmelo/ecoaplic/r
2 penguins_df
```

	X	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm
1	1	Adelie	Torgersen	39.1	18.7	181
2	2	Adelie	Torgersen	39.5	17.4	186
3	3	Adelie	Torgersen	40.3	18.0	195
4	4	Adelie	Torgersen	NA	NA	NA
5	5	Adelie	Torgersen	36.7	19.3	193
6	6	Adelie	Torgersen	39.3	20.6	190
7	7	Adelie	Torgersen	38.9	17.8	181
8	8	Adelie	Torgersen	39.2	19.6	195
9	9	Adelie	Torgersen	34.1	18.1	193
10	10	Adelie	Torgersen	42.0	20.2	190
11	11	Adelie	Torgersen	37.8	17.1	186
12	12	Adelie	Torgersen	37.8	17.3	180
13	13	Adelie	Torgersen	41.1	17.6	182
14	14	Adelie	Torgersen	38.6	21.2	191
15	15	Adelie	Torgersen	34.6	21.1	188

Planning a data visualization



source: Andrew Gard

<https://www.youtube.com/@EquitableEquations>

We know this data

```
1 library(tidyverse)
2 library(palmerpenguins)
3
4 data("penguins")
5 penguins %>%
6   select(1:5)
```

A tibble: 344 × 5

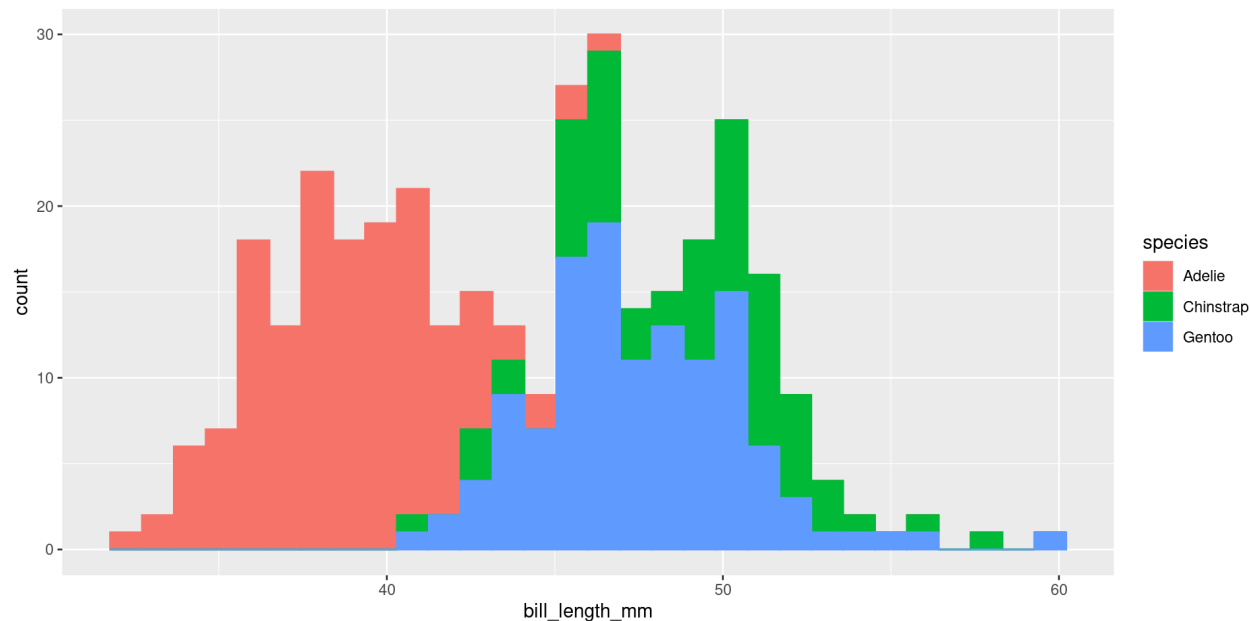
	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm
	<fct>	<fct>	<dbl>	<dbl>	<int>
1	Adelie	Torgersen	39.1	18.7	181
2	Adelie	Torgersen	39.5	17.4	186
3	Adelie	Torgersen	40.3	18	195
4	Adelie	Torgersen	NA	NA	NA
5	Adelie	Torgersen	36.7	19.3	193
6	Adelie	Torgersen	39.3	20.6	190
7	Adelie	Torgersen	38.9	17.8	181
8	Adelie	Torgersen	39.2	19.6	195
9	Adelie	Torgersen	34.1	18.1	193
10	Adelie	Torgersen	42	20.2	190

i 334 more rows

How to visually check continuous variables?

Histograms

```
1 library(tidyverse)
2 library(palmerpenguins)
3
4 data("penguins")
5 penguins %>%
6   group_by(species) %>%
7     ggplot(aes(x=bill_length_mm, color=species, fill=species))+
8     geom_histogram()
```

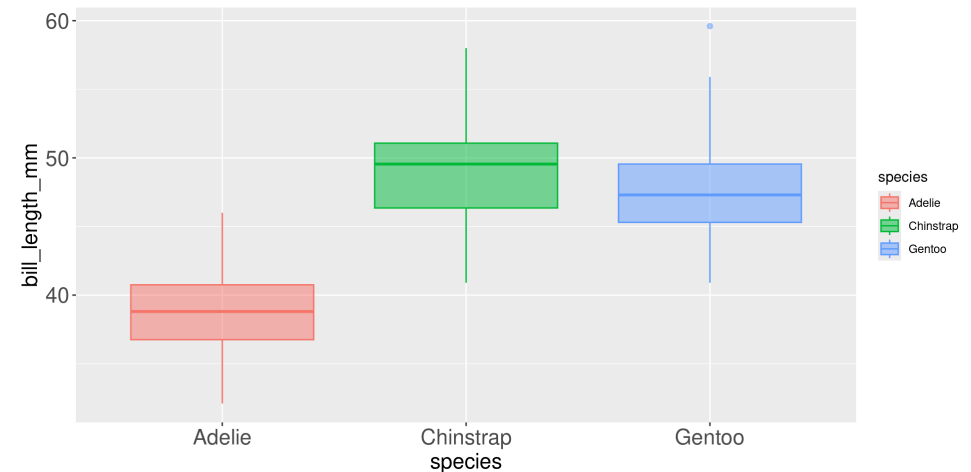


Boxplots

```

1 library(tidyverse)
2 library(palmerpenguins)
3
4 data("penguins")
5 penguins %>%
6   group_by(species) %>%
7     ggplot(aes(x=species,
8               y=bill_length_mm,
9               color=species,
10              fill=species))+
11     geom_boxplot(alpha=0.5)+
12     theme(axis.text=element_text(size=12),
13           axis.title=element_text(size=14))

```



Your turn

- Try to reproduce with any other continuous variable
- |Do a Histogram and a Boxplot

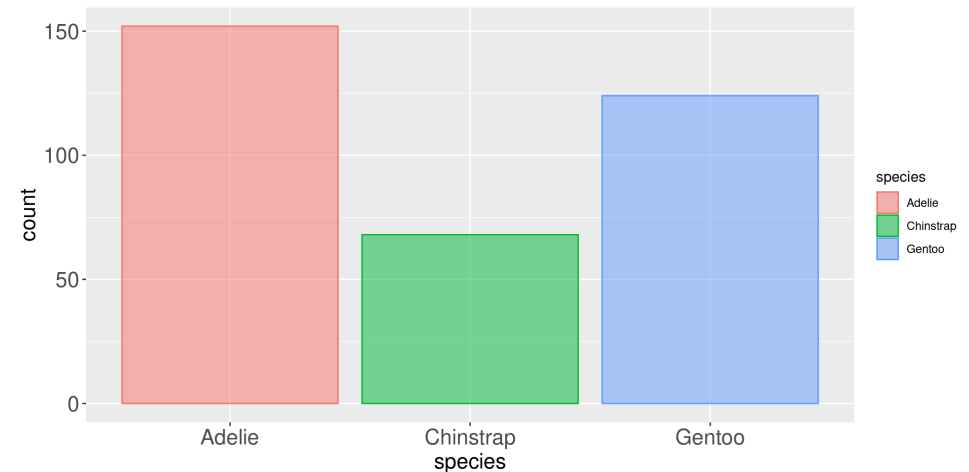
Checking categorical variables

Species of penguin

```

1 library(tidyverse)
2 library(palmerpenguins)
3
4 penguins %>%
5   ggplot(aes(x=species,
6             color=species,
7             fill=species))+
8   geom_bar(alpha=0.5)+
9   theme(axis.text=element_text(size=12),
10         axis.title=element_text(size=14))

```



Observations per year

```

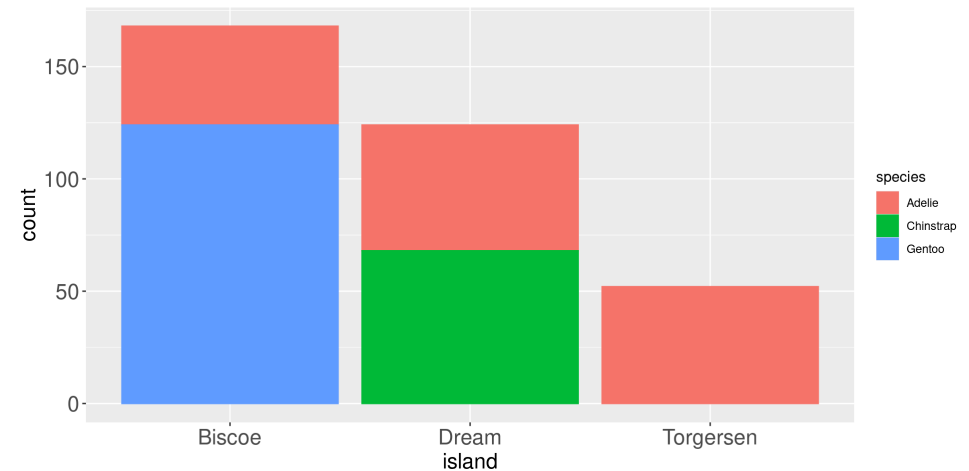
1 library(tidyverse)
2 library(palmerpenguins)
3
4 penguins %>%
5   ggplot(aes(x=year,
6             color=species,
7             fill=species))+
8   geom_bar()+
9   theme(axis.text=element_text(size=12),
10         axis.title=element_text(size=14))

```



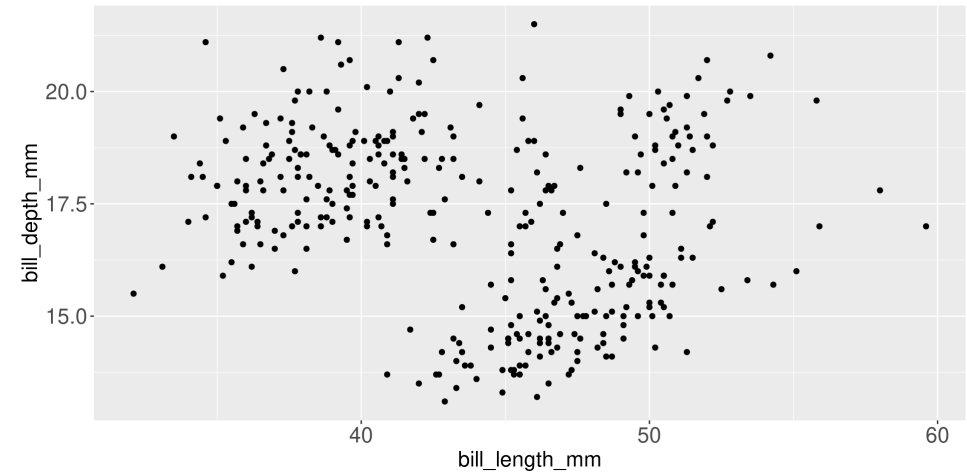
Observations per island

```
1 library(tidyverse)
2 library(palmerpenguins)
3
4 penguins %>%
5   ggplot(aes(x=island,
6             color=species,
7             fill=species))+
8   geom_bar()+
9   theme(axis.text=element_text(size=12),
10         axis.title=element_text(size=14))
```



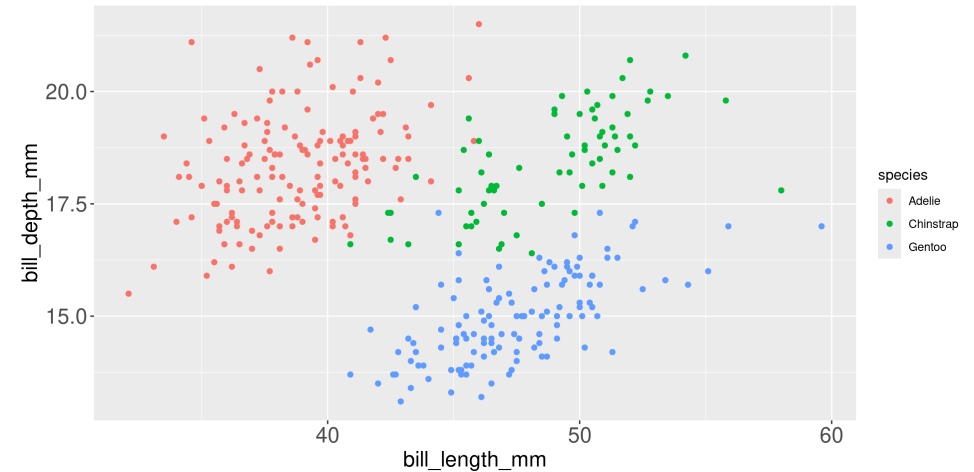
Visualising correlations

```
1 penguins %>%  
2   ggplot(aes(x=bill_length_mm,  
3             y = bill_depth_mm))+  
4   geom_point()+  
5   theme(axis.text=element_text(siz  
6         axis.title=element_text(si
```



Visualising correlations per species

```
1 penguins %>%  
2   ggplot(aes(x=bill_length_mm,  
3             y = bill_depth_mm,  
4             color=species,  
5             fill=species))+  
6   geom_point()+  
7   theme(axis.text=element_text(siz  
8         axis.title=element_text(si
```

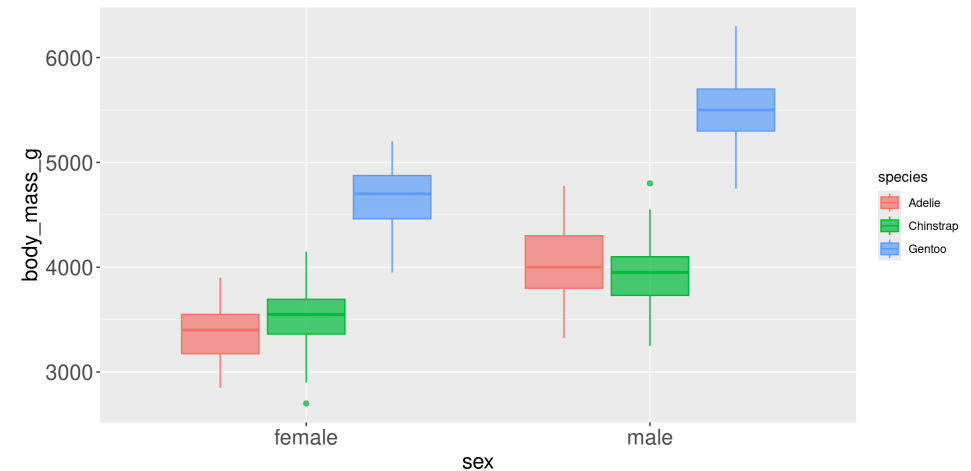


Body mass per sex

```

1 penguins %>%
2   na.omit() %>%
3   ggplot(aes(x=sex,
4             y = body_mass_g,
5             color=species,
6             fill=species))+
7   geom_boxplot(alpha=0.7)+
8   theme(axis.text=element_text(size=12),
9         axis.title=element_text(size=14))

```

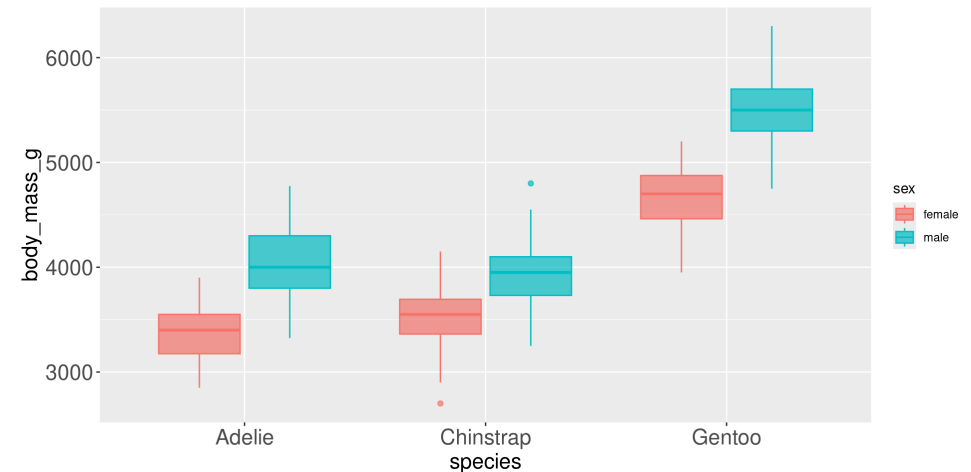


Body mass per sex (iverting groups)

```

1 penguins %>%
2   na.omit() %>%
3   ggplot(aes(x=species,
4             y = body_mass_g,
5             color=sex,
6             fill=sex))+
7   geom_boxplot(alpha=0.7)+
8   theme(axis.text=element_text(size=12),
9         axis.title=element_text(size=14))

```



Your turn

- Can body mass predict bill length?
- Do sex explain flipper length

Exploring data is asking relevant questions

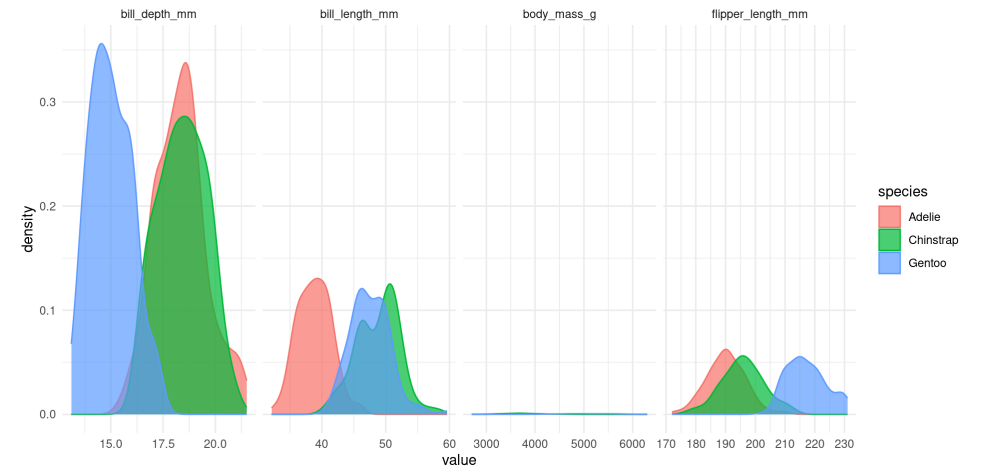
- This is not mining
- Don't just correlate random things
- Start to imagine before coding

Check distributions

```

1 penguins %>%
2   na.omit() %>%
3   pivot_longer(bill_length_mm:body
4   ggplot(aes(x=value,
5             group=species,
6             fill=species,
7             color=species))+
8   geom_density(alpha=0.7)+
9   facet_grid(~trait, scales = "fre
10  theme(axis.text=element_text(siz
11        axis.title=element_text(si
12  theme_minimal()

```



The importance of distributions

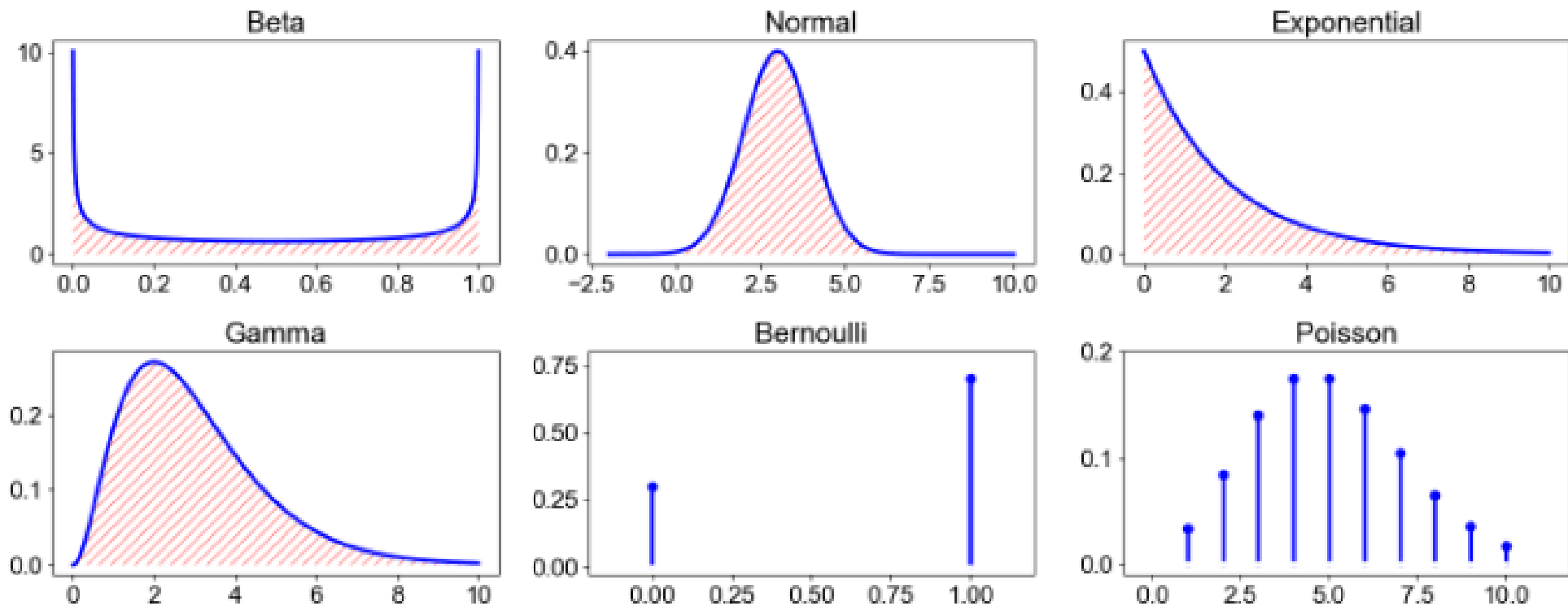
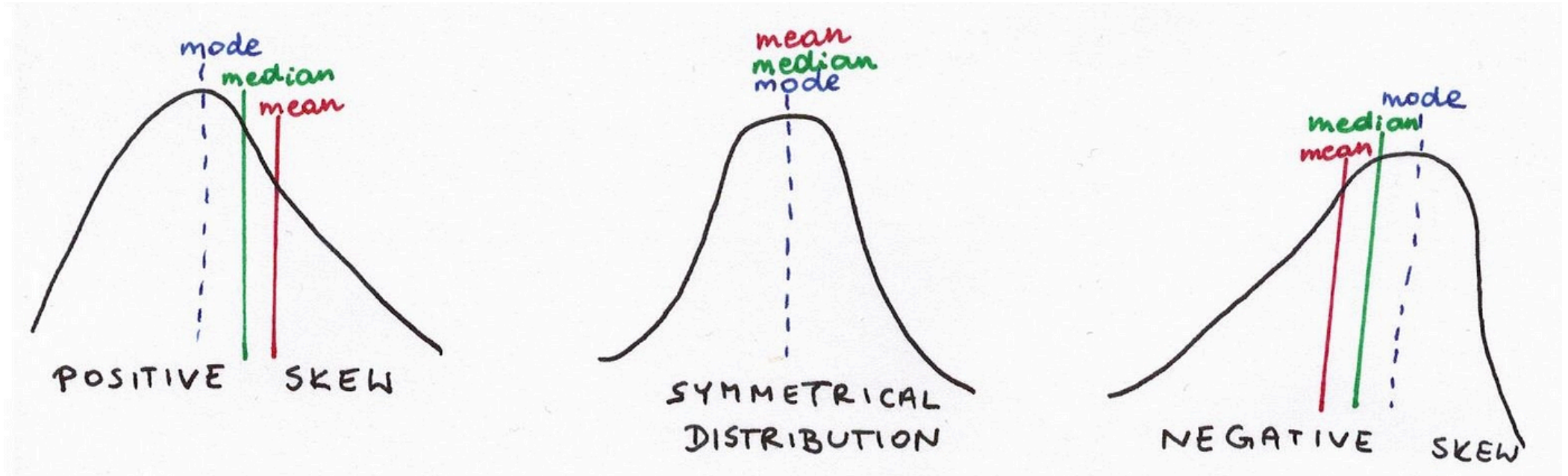


Figure 1. Beta, normal, exponential, gamma, Bernoulli, and Poisson distributions, each with a total mass of one.

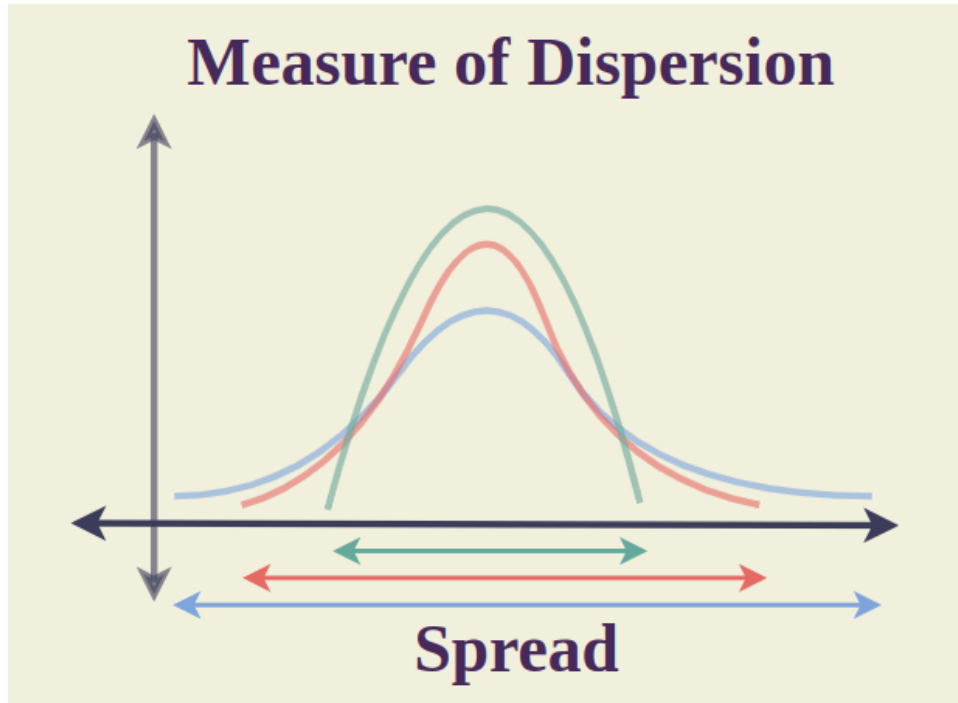
source: <https://gregorygundersen.com/blog/2020/04/11/moments/>

Moments of centrality



Mean, median and mode

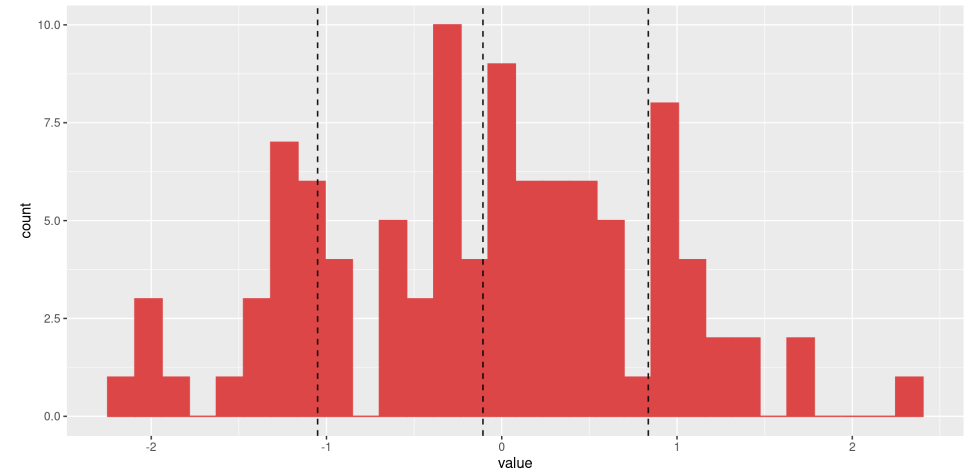
Moments of dispersion



- Variance
- Standard deviation
- Standard Error
- Range
- Quantiles

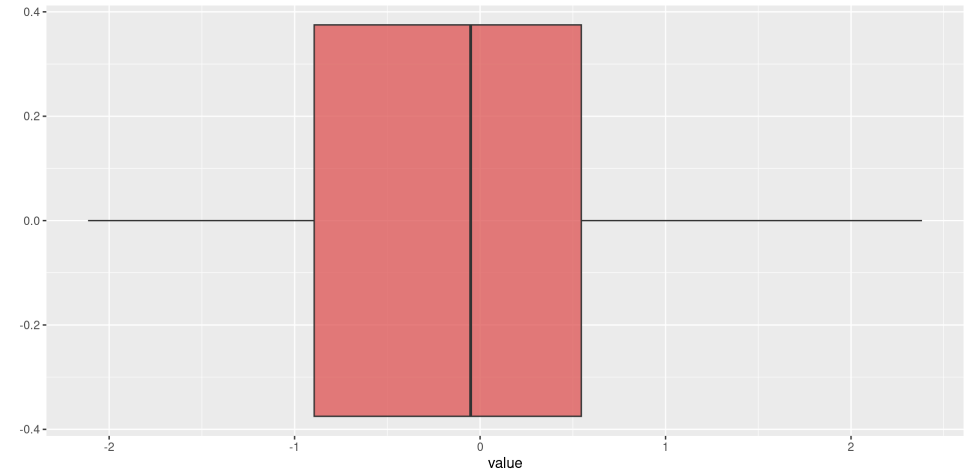
Checking via histogram

```
1 set.seed(999)
2 normal<-rnorm(100)
3 normal %>%
4   as.tibble() %>%
5   ggplot(aes(value))+
6   geom_histogram(color="#DD4A48",
7   geom_vline(xintercept=c(mean(normal),
8               linetype="dashed"))
```

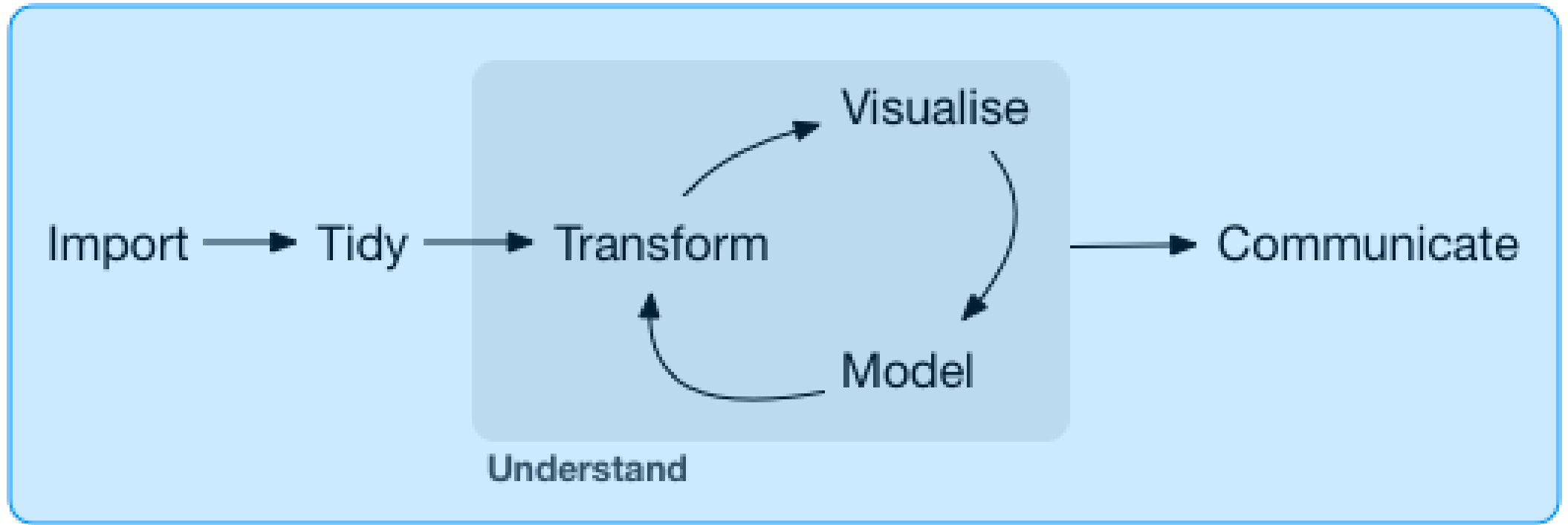


Checking via bokplot

```
1 set.seed(999)
2 normal<-rnorm(100)
3 normal %>%
4   as.tibble() %>%
5   ggplot(aes(value))+
6   geom_boxplot(fill="#DD4A48", alph
```



Workflow



Program

End of session on DA