

VISUALISATION IN R (FOCUS ON GGPLOTS)

Rstats Friday! | Dr. Erin I. Walsh

Session overview

- The thoughtful figure

- Philosophy

Key source: Tufte, E., & Graves-Morris, P. (2014). *The visual display of quantitative information.*; 1983

- Base R graphics

- The grammar of graphics

(aka How I learned to stop worrying and learned to love ggplots)

- Grammar

Key source: Wilkinson, L. (2006). *The grammar of graphics*. Springer Science & Business Media. (or the 2010 paper)

- Ggplot2

Key source: Elegant graphics for data analysis

- Q & A

- Optional bonus slides on the grammar of graphics if you're gluttons for punishment

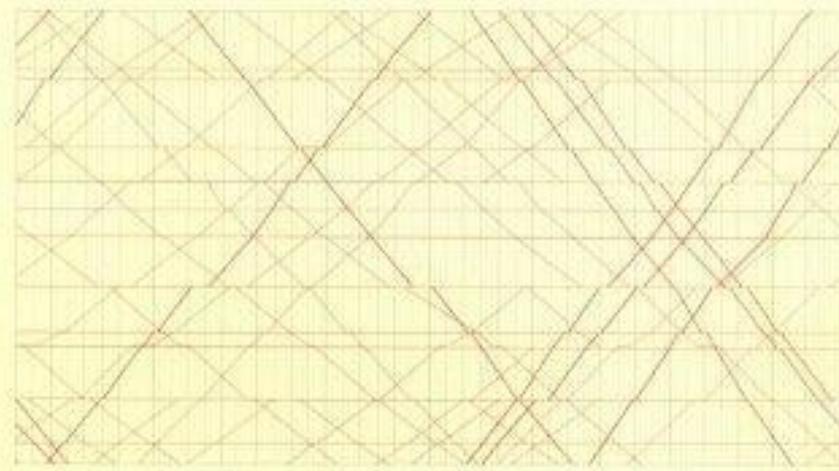
Graphics aren't just for Christmas

- Foundation of your data exploration/analysis
- Part of communicating with colleagues and yourself (record keeping) as you work on a project
- Should be drafted like any other part of the manuscript
 - Iteratively
 - Not an afterthought!



Tufte, E., &
Graves-Morris,
P. (2014). The
visual display of
quantitative
information.;
1983

Copyrighted Material



SECOND EDITION

The Visual Display
of Quantitative Information

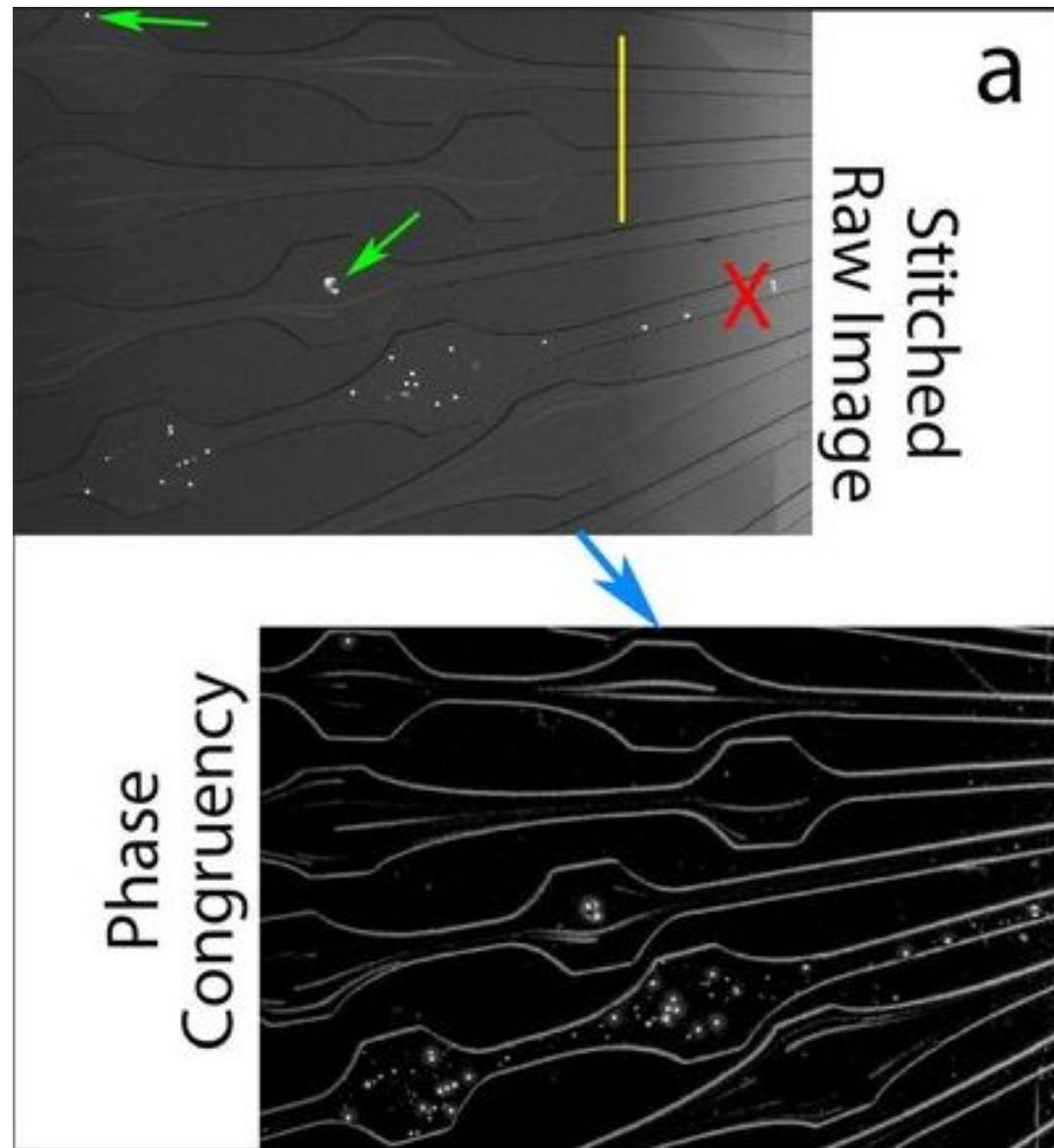
EDWARD R. TUFTE

Copyrighted Material

Principles of Graphical Excellence

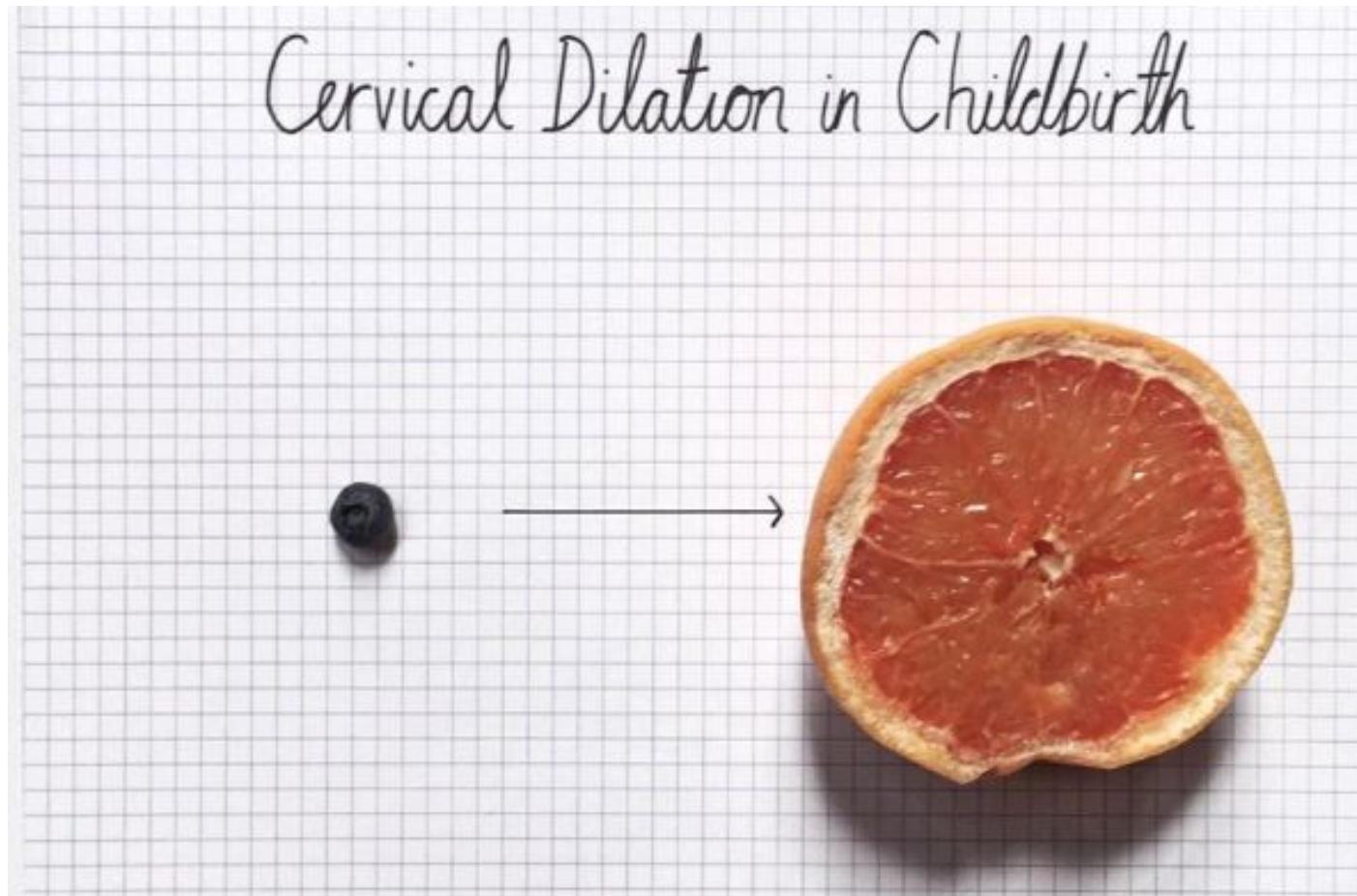
- Show the data
- Induce the viewer to think about the substance rather than about methodology, graphic design, the technology of the graphic production or something else
- Avoid distorting what the data have to say
- Present many numbers in a small space
- Make large datasets coherent
- Encourage the eye to compare different pieces of data
- Reveal the data at several levels of detail, from a broad overview to a fine structure
- Serve a reasonably clear purpose: description, exploration, tabulation or decoration
- Be closely integrated with the statistical verbal descriptions of a dataset

Show the data



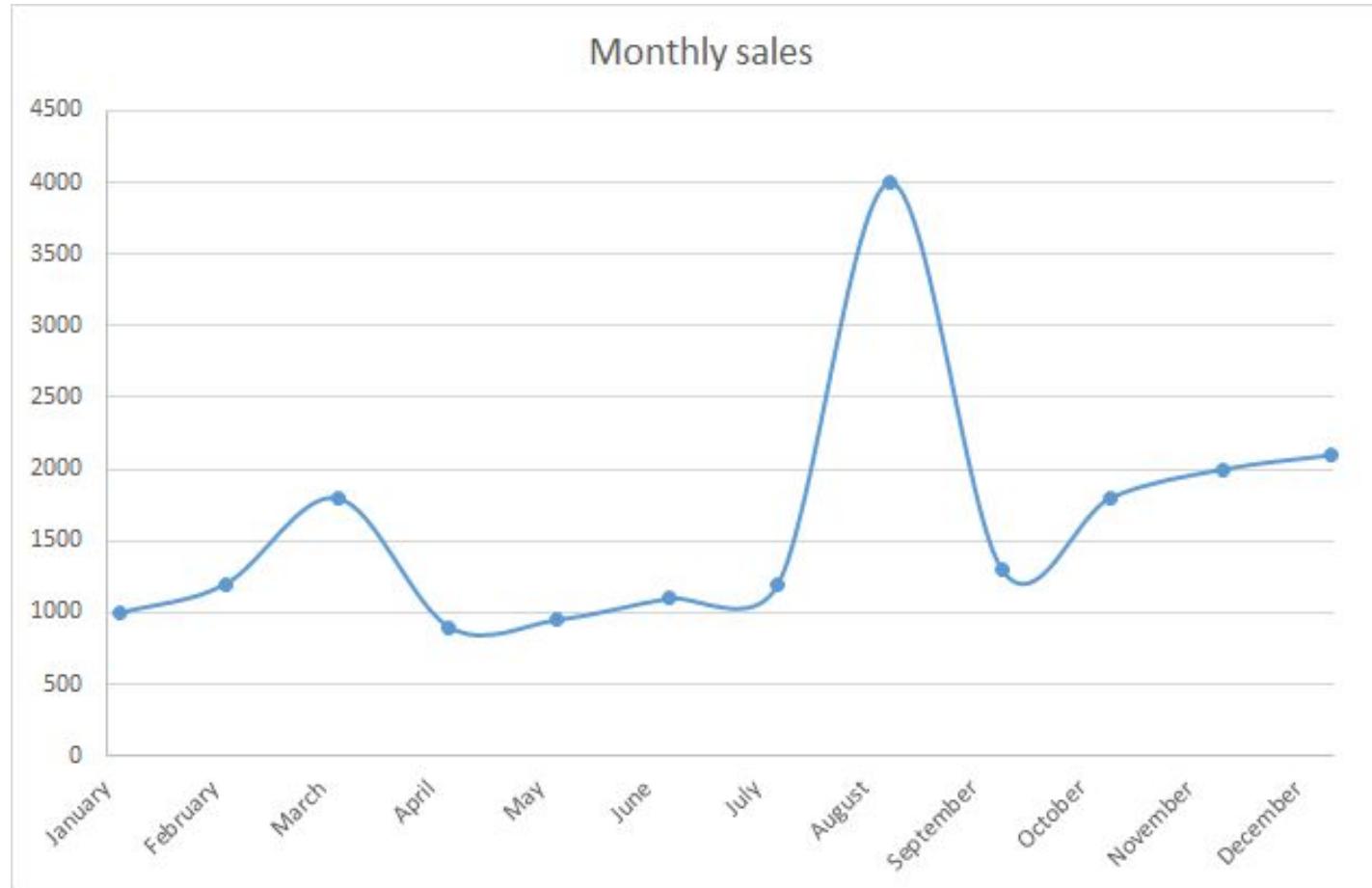
Example from: Benson, B. L., Li, L., Myers, J. T., Dorand, R. D., Gurkan, U. A., Huang, A. Y., & Ransohoff, R. M. (2018). Biomimetic post-capillary venule expansions for leukocyte adhesion studies. *Scientific Reports*, 8(1), 9328.

Induce the viewer to think about the substance



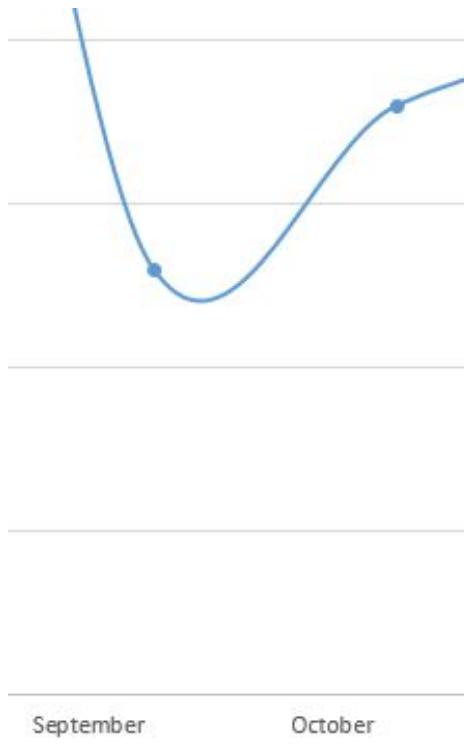
Example from: Mona Chalabi from
<http://sites.uci.edu/fmd2017fa/files/2017/05/Illustrator-2.pdf>

Avoid distorting what the data have to say



Example from: <https://m.signalvnoise.com/lets-chart-stop-those-lying-line-charts/>

Avoid distorting what the data have to say



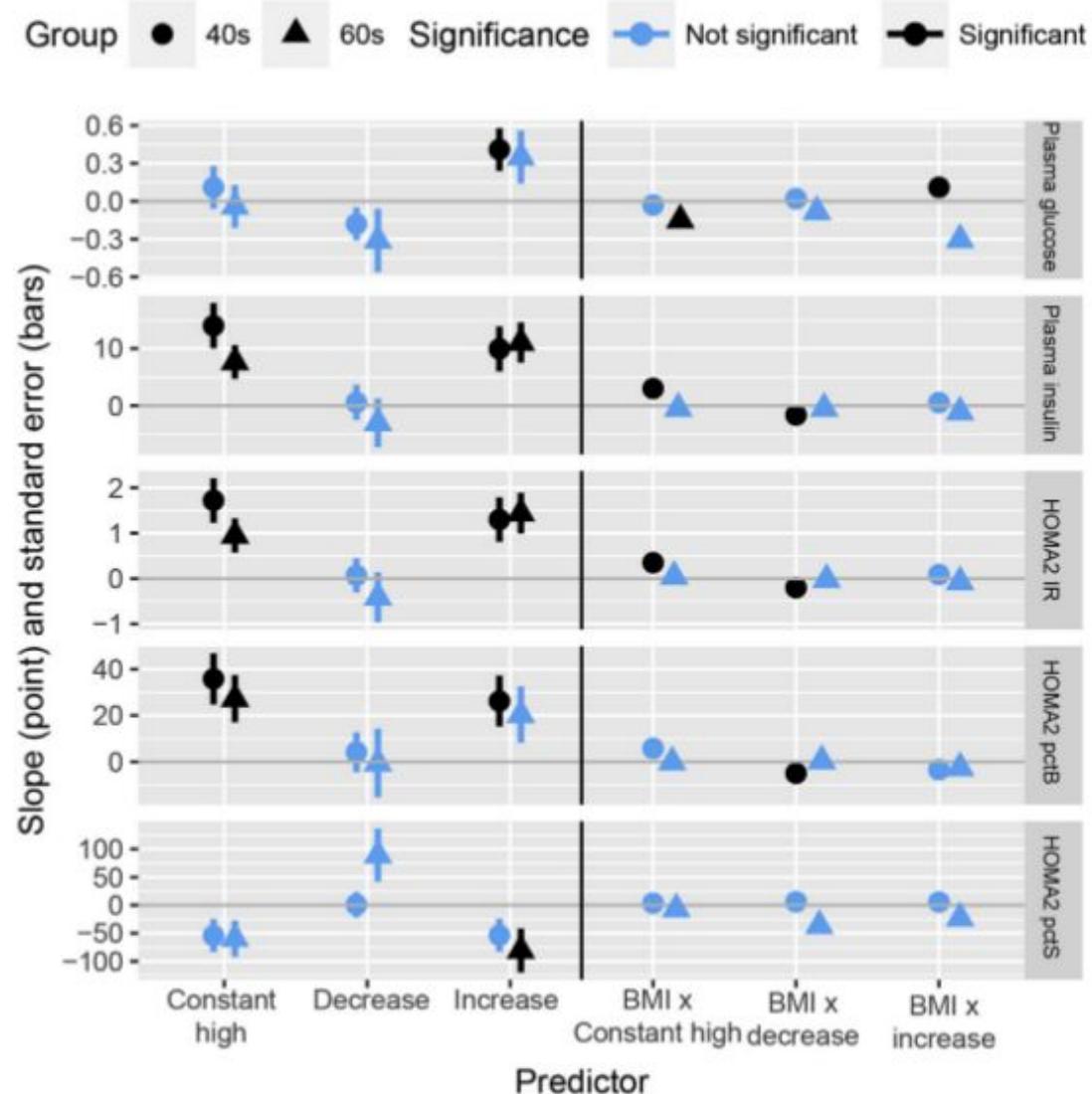
Example from: <https://m.signalvnoise.com/lets-chart-stop-those-lying-line-charts/>

Present many numbers in a small space

Table 3. Individual and cohort analysis.

Outcome	Predictor	40s			60s			
		b	SE	p	b	SE	p	
Model (1)								
Plasma glucose	Intercept	4.20	0.18	<0.01	4.74	0.31	<0.01	
	HOMA2 IR	0.03	0.01	0.16	0.02	0.01	0.02	
	Adj R ²	0.36			0.14			
Plasma insulin	Intercept	12.3	0.23	<0.01	11.11	0.42	<0.01	
	HOMA2 IR	0.11	0.01	<0.01	1.01	0.14	<0.01	
	Adj R ²	0.27			0.23			
HOMA2 IR	Intercept	1.34	0.03	<0.01	14.1	1.08	<0.01	
	HOMA2 IR	0.27	0.02	<0.01	0.11	0.02	<0.01	
	Adj R ²	0.20			0.21			
HOMA2 pctB	Intercept	1.34	0.03	<0.01	11.08	0.37	<0.01	
	HOMA2 pctB	0.04	0.02	0.18	0.01	0.01	0.01	
	Adj R ²	0.24			0.18			
HOMA2 pctB	Intercept	1.30-33	0.03	<0.01	22.21	0.46	<0.01	
	HOMA2 pctB	0.03	0.01	<0.01	0.06	0.01	<0.01	
	Adj R ²	0.22			0.21			
Model (2)	Plasma glucose	Intercept	0.01	0.01	<0.01	4.90	0.31	<0.01
	Constant HOMA2 IR	0.11	0.01	<0.01	0.12	0.17	0.01	
	HOMA2 IR	0.19	0.03	<0.01	0.18	0.20	0.01	
	Intercept	0.41	0.03	<0.01	0.15	0.21	0.01	
Plasma insulin	Intercept	11.03	0.13	<0.01	10.35	0.36	<0.01	
	Constant HOMA2 IR	1.01	0.08	<0.01	1.01	2.08	0.01	
	HOMA2 IR	0.05	0.02	<0.01	0.05	0.08	0.01	
	Intercept	0.03	0.01	<0.01	11.06	1.03	<0.01	
HOMA2 IR	Intercept	1.04	0.14	<0.01	2.4	1.36	<0.01	
	HOMA2 IR	0.22	0.03	<0.01	0.01	0.38	0.01	
	Intercept	1.2	0.05	<0.01	14.4	0.45	<0.01	
HOMA2 pctB	Intercept	0.1	0.01	<0.01	0.1	0.01	<0.01	
	HOMA2 pctB	1.26-36	20.3	<0.01	11.11	1.14	<0.01	
	Decrease	0.01	0.01	<0.01	21.02	1.07	<0.01	
	Intercept	0.22	0.01	<0.01	0.01	0.01	<0.01	
	HOMA2 pctB	20.19	11.00	<0.01	20.31	12.22	<0.01	
HOMA2 pctB	Intercept	11.73	0.03	<0.01	-0.03	0.01	<0.01	
	Constant HOMA2 IR	0.01	0.01	<0.01	0.01	0.01	<0.01	
	HOMA2 IR	0.01	0.01	<0.01	0.01	0.01	<0.01	
	Intercept	0.01	0.01	<0.01	0.01	0.01	<0.01	
HOMA2 pctB	Intercept	0.01	0.01	<0.01	0.01	0.01	<0.01	
	HOMA2 pctB	0.01	0.01	<0.01	0.01	0.01	<0.01	
Model (3)	Plasma glucose	Intercept	4.06	0.18	<0.01	3.21	0.03	<0.01
	HOMA2 IR	0.01	0.01	0.09	0.12	0.06	0.01	
	HOMA2 pctB	0.02	0.01	0.01	0.02	0.06	0.01	
	Adj R ²	0.11	0.04	0.01	0.04	0.06	0.01	
Plasma insulin	Intercept	14.08	0.10	<0.01	11.46	0.31	<0.01	
	HOMA2 IR	0.02	0.01	0.02	0.07	0.04	0.01	
	HOMA2 pctB	0.02	0.01	0.02	0.01	0.04	0.01	
	Adj R ²	0.12	0.04	0.01	0.06	0.06	0.01	
HOMA2 IR	Intercept	1.81	0.01	<0.01	14.0	0.06	0.2	
	HOMA2 IR	0.01	0.01	0.01	0.00	0.01	0.01	
	HOMA2 pctB	0.02	0.01	0.02	0.01	0.01	0.01	
	Adj R ²	0.13	0.04	0.01	0.07	0.01	0.01	
HOMA2 pctB	Intercept	1.20-28	21.44	<0.01	10.33	3.30	<0.01	
	HOMA2 pctB	0.01	0.01	0.01	0.01	0.01	0.01	
	HOMA2 pctB	0.01	0.01	0.01	0.01	0.01	0.01	
	Adj R ²	0.14	0.04	0.01	0.08	0.01	0.01	
HOMA2 pctB	Intercept	1.20-33	21.44	<0.01	10.33	3.30	<0.01	
	HOMA2 pctB	0.01	0.01	0.01	0.01	0.01	0.01	
	Adj R ²	0.14	0.04	0.01	0.08	0.01	0.01	
HOMA2 pctB	Intercept	1.20-33	21.44	<0.01	10.33	3.30	<0.01	
	HOMA2 pctB	0.01	0.01	0.01	0.01	0.01	0.01	
	Adj R ²	0.14	0.04	0.01	0.08	0.01	0.01	
Legend: PA = Physical activity; HOMA2 IR = Homeostatic model assessment HOMA2 IR; HOMA2 pctB and HOMA2 pcts. These models were fit to the complete dataset. All models adjust for age, gender, log-transformed pre-diabetes status, four group comparisons in overall fasting glucose (HOMA2 = required fasting glucose) and physical activity.								

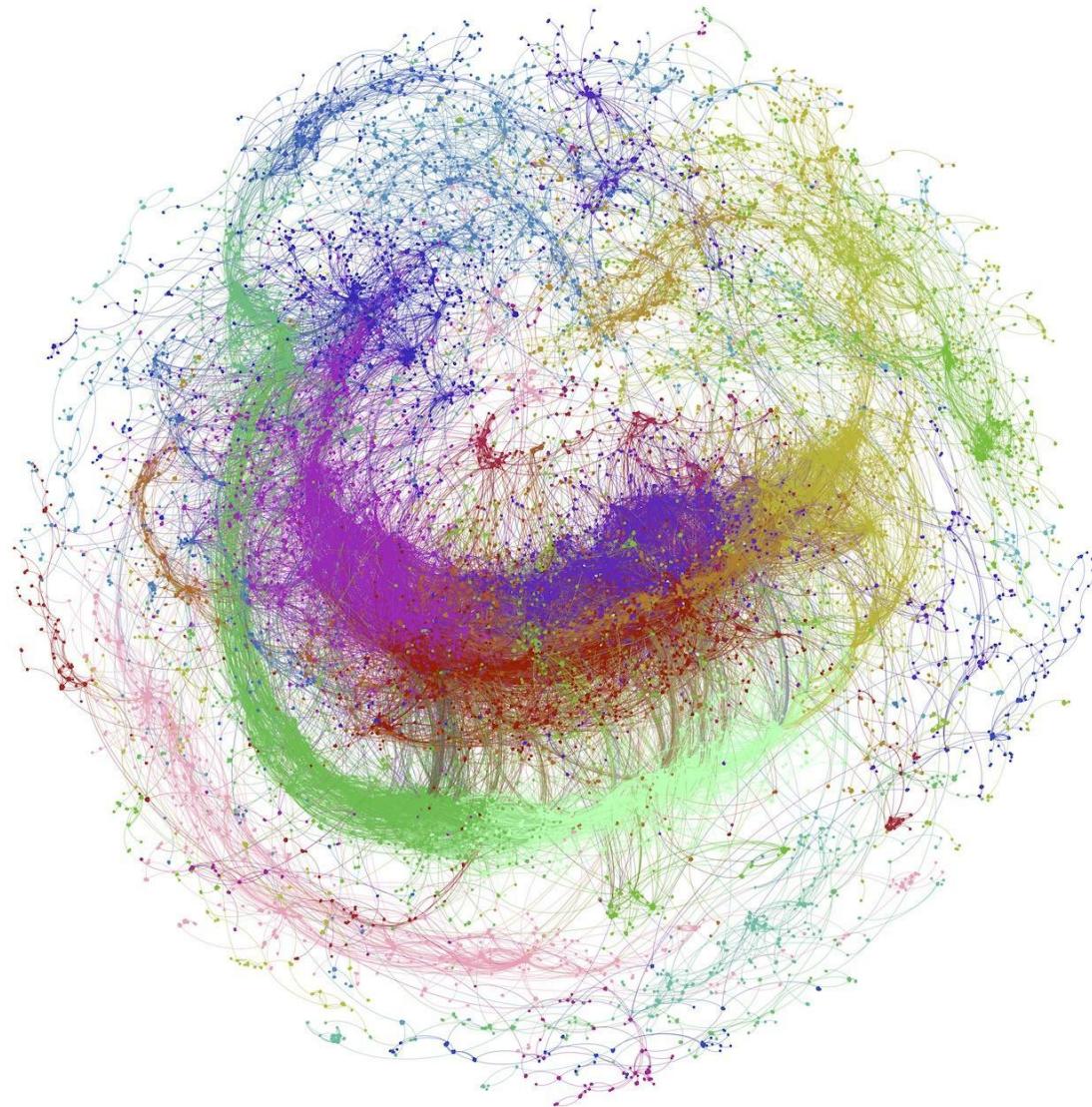
Example from: Walsh, E. I., Shaw, J., & Cherbuin, N. (2018). Trajectories of BMI change impact glucose and insulin metabolism. *Nutrition, Metabolism and Cardiovascular Diseases*, 28(3), 243-251.



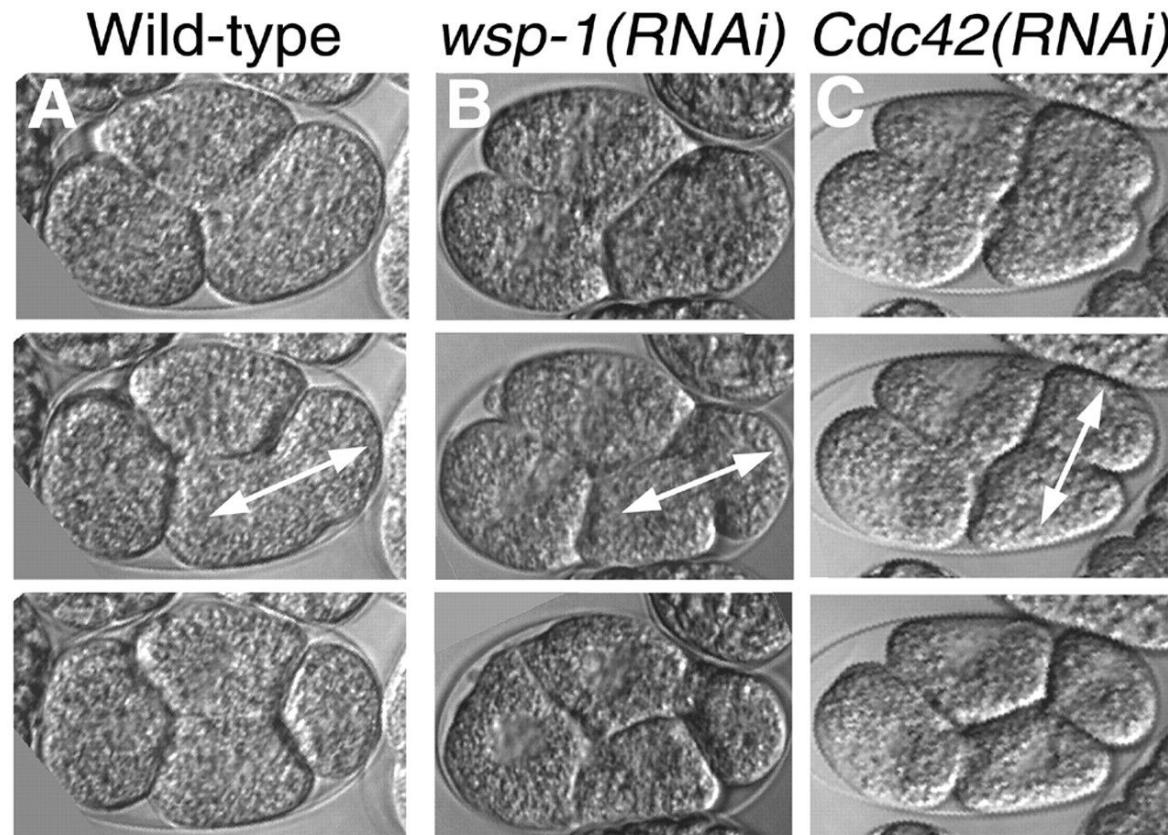
Make large datasets coherent

Example from: GBR coherent patterns from Higgs dataset
<https://lts2.epfl.ch/blog/gbr/2015/05/01/principal-patterns-on-graphs-discovering-coherent-structures-in-datasets/>

Benzi, K., Ricaud, B., & Vanderghenst, P. (2016). Principal Patterns on Graphs: Discovering Coherent Structures in Datasets. *IEEE Trans. Signal and Information Processing over Networks*, 2(2), 160-173.

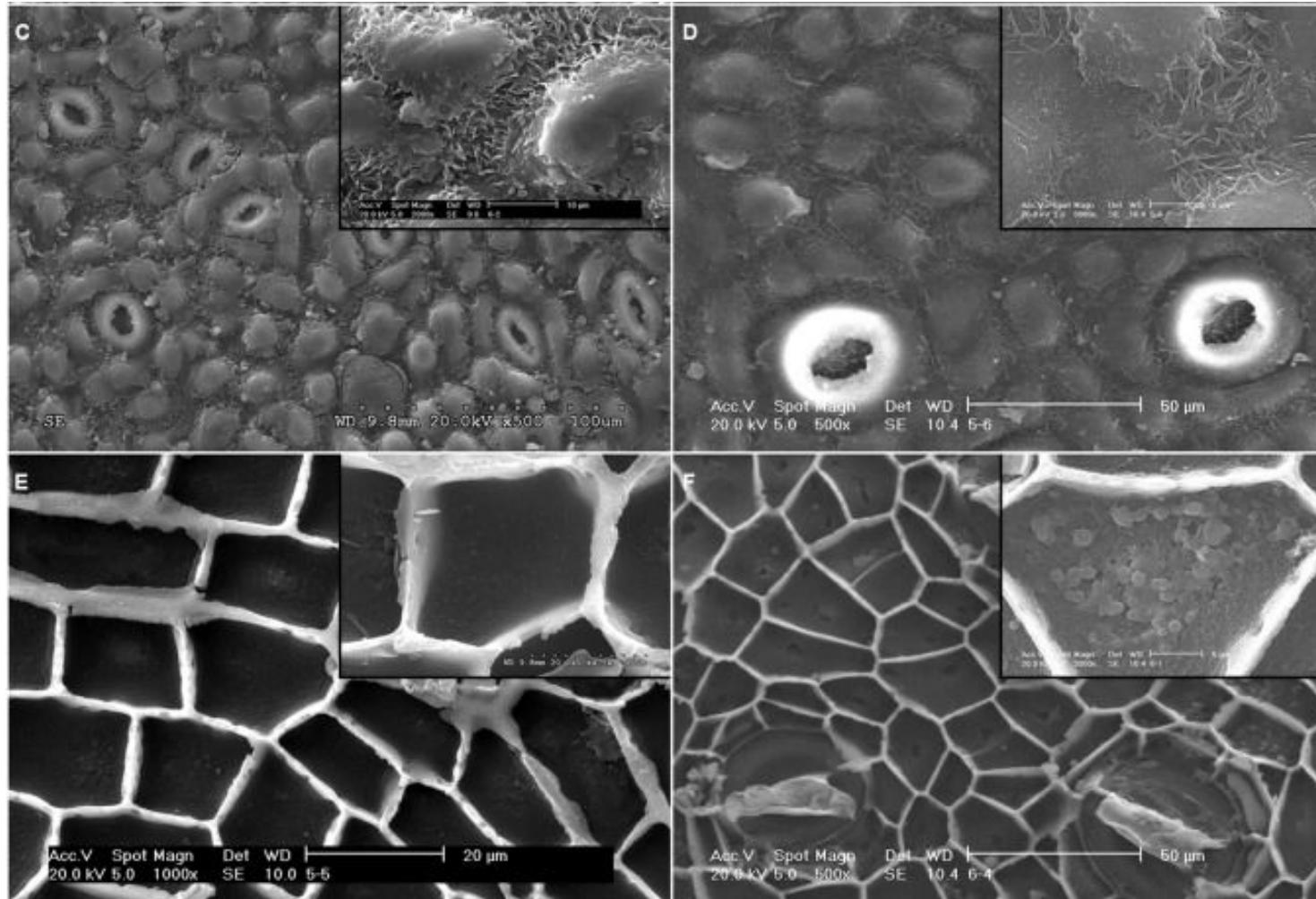


Encourage the eye to compare different pieces of data



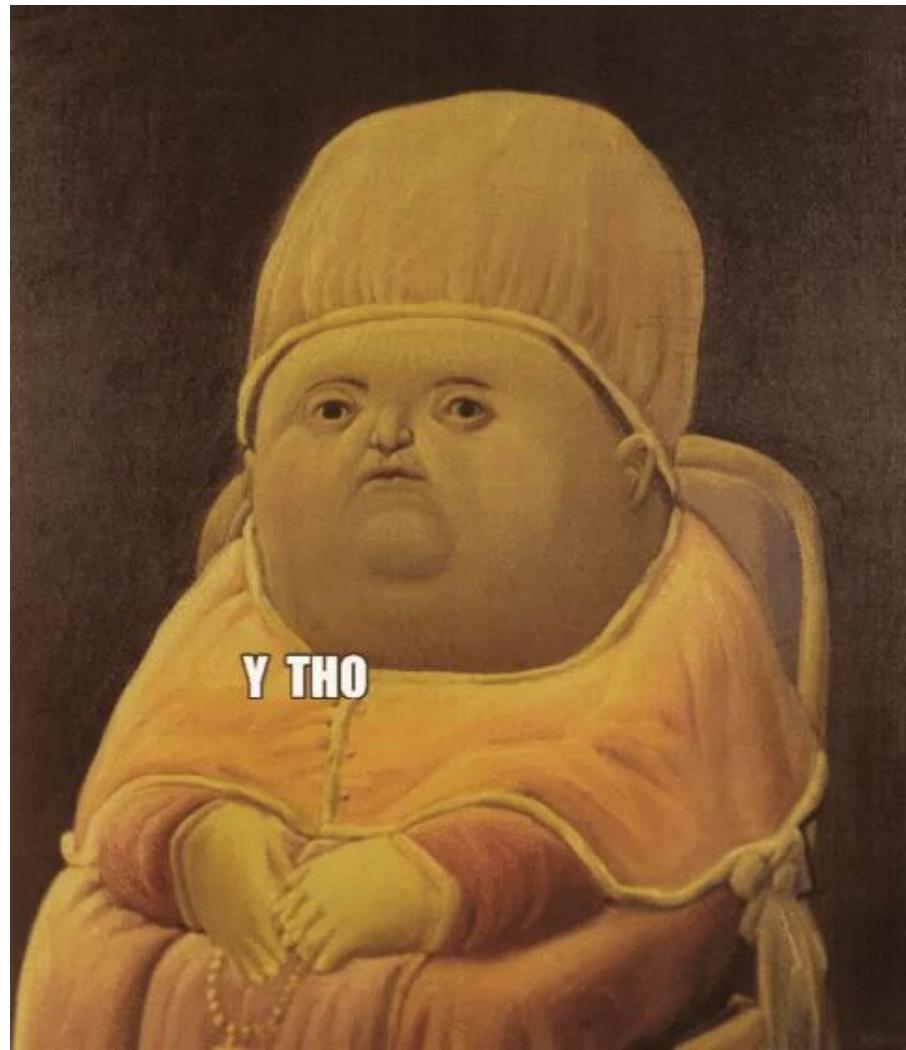
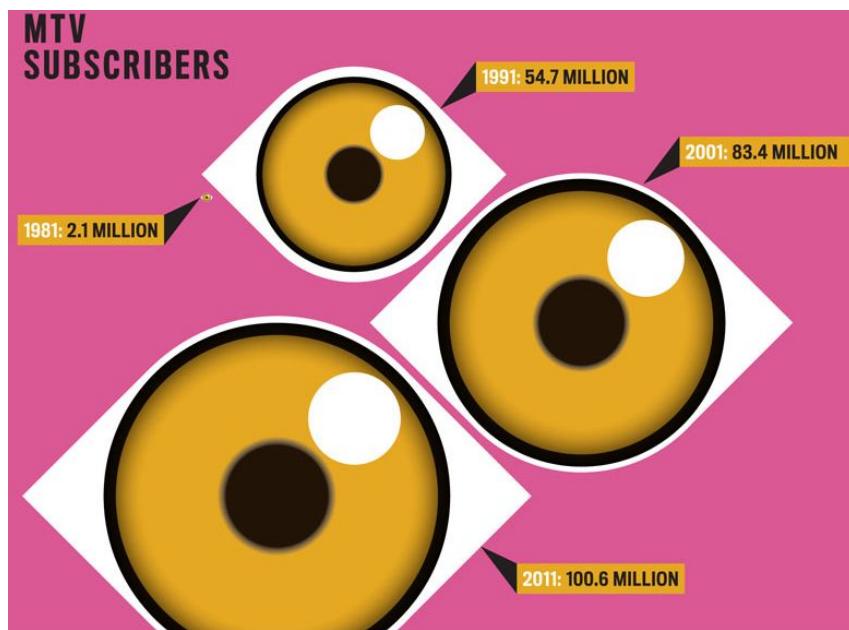
Example from: Sawa, M., Suetsugu, S., Sugimoto, A., Miki, H., Yamamoto, M., & Takenawa, T. (2003). Essential role of the *C. elegans* Arp2/3 complex in cell migration during ventral enclosure. *Journal of cell science*, 116(8), 1505-1518. <http://jcs.biologists.org/content/116/8/1505>

Reveal the data at several levels of detail



Example from:
Guzmán, P.,
Fernández, V.,
Graça, J., Cabral,
V., Kayali, N.,
Khayet, M., & Gil,
L. (2014). Chemical
and structural
analysis of
Eucalyptus
globulus and *E.*
camaldulensis leaf
cuticles: a lipidized
cell wall region.
Frontiers in plant
science, 5, 481.

Serve a reasonably clear purpose

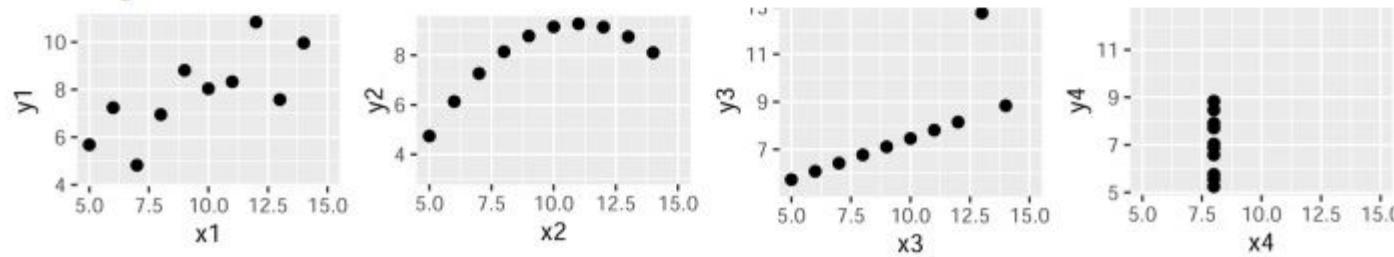


Be closely integrated with the statistical verbal descriptions of a dataset

Viewed as pure numbers, it is difficult to see any difference between the sets:

x1	y1	x2	y2	x3	y3	x4	y4
10	8.04	10	9.14	10	7.46	8	6.58
8	6.95	8	8.14	8	6.77	8	5.76
13	7.58	13	8.74	13	12.74	8	7.71
9	8.81	9	8.77	9	7.11	8	8.84
11	8.33	11	9.26	11	7.81	8	8.47
14	9.96	14	8.1	14	8.84	8	7.04
6	7.24	6	6.13	6	6.08	8	5.25
4	4.26	4	3.1	4	5.39	19	12.5
12	10.84	12	9.13	12	8.15	8	5.56
7	4.82	7	7.26	7	6.42	8	7.91
5	5.68	5	4.74	5	5.73	8	6.89

Striking differences become immediately obvious once they are displayed as scatterplots.



Base R Graphics

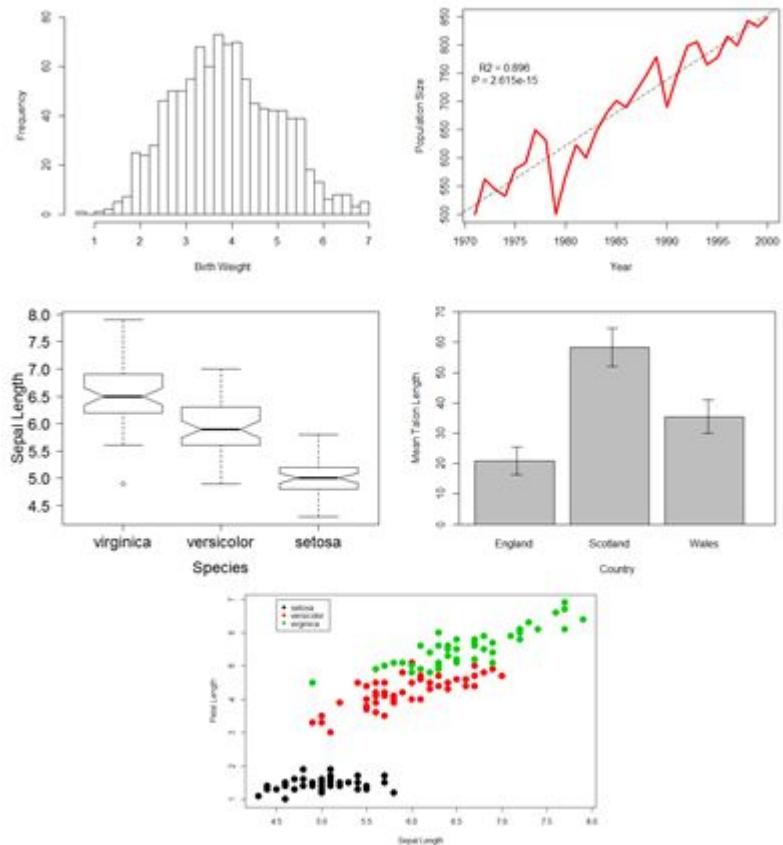
Pros:

- Quick and simple
- Most basic needs implemented
- Extremely convenient

Cons:

- Unusual plots difficult to create
- Not going to win any pageants
 - (but good if you get stuck on aesthetics when you should be focussing on the next step of analysis).

Recommended usage scenario: Part of initial data exploration and discovery.



Exercise 01 time!

Single variable plots

- Histograms
- Dot plots
- Strip charts

Multiple variable plots

- Scatterplots
- Paired scatterplots
- Boxplots

Output of other functions or operations

- Density plots
- Linear model diagnostics

Bar plots

Load the example data

```
data("iris")
```

```
iris$Species<-as.factor(iris$Species)
```

```
head(iris)
```

```
> head(iris)
   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
1          5.1         3.5          1.4         0.2  setosa
2          4.9         3.0          1.4         0.2  setosa
3          4.7         3.2          1.3         0.2  setosa
4          4.6         3.1          1.5         0.2  setosa
5          5.0         3.6          1.4         0.2  setosa
6          5.4         3.9          1.7         0.4  setosa
```

Single variable plots

Histogram

```
hist(iris$Sepal.Length)
```

Dot chart

```
dotchart(iris$Sepal.Length)
```

Strip chart

```
stripchart(iris$Sepal.Length)
```

Multiple variable plots

Scatterplot

```
plot(iris$Petal.Length, iris$Petal.Width)
```

Adding lines to a scatterplot

```
plot.new()  
lines(iris$Petal.Length, iris$Petal.Width)  
plot(iris$Petal.Length, iris$Petal.Width)  
lines(iris$Petal.Length, iris$Petal.Width)  
  
plot(iris$Petal.Length, iris$Petal.Width, type="l")
```

Multiple variable plots

Pairs plot

```
pairs(~iris$Sepal.Length +  
      iris$Petal.Width +  
      iris$Sepal.Width)
```

Box plot

```
plot(iris$Species, iris$Sepal.Length, type="b")  
plot(iris$Species, iris$Sepal.Length, type="b",  
     notch=TRUE)
```

Output of other functions or operations

Density plot (cousin of the histogram)

```
sepal_density<-density(iris$Sepal.Length)
plot(sepal_density)
```

Linear model diagnostics

```
sepal_linear<-lm(iris$Sepal.Length ~iris$Petal.Width)
plot(sepal_linear)
```

Linear model output on points

```
plot(iris$Sepal.Length ~iris$Petal.Width)
abline(sepal_linear)
```

Output of other functions or operations

Bar plots

```
sepal_table<-table(iris$Sepal.Length)
barplot(sepal_table)
barplot(sepal_table,horiz=TRUE)
```

Stacked bar plots

```
sepal_table_stacked<-table(iris$Species,
                           iris$Sepal.Length)
barplot(sepal_table_stacked)
barplot(sepal_table_stacked,
        legend = names(table(iris$Species)))
```

Pie chart

```
petal_width_bins<-rep("Under 1", length(iris$Petal.Width))
petal_width_bins[which(iris$Petal.Width>1)]<-"Over 1"
petal_width_bins[which(iris$Petal.Width>2)]<-"Over 2"
pie(table(petal_width_bins))
```

Multiple plots

```
par(mfrow=c(1,1))
hist(iris$Sepal.Length)

par(mfrow=c(1,2))
hist(iris$Sepal.Length)
hist(iris$Sepal.Length)

par(mfrow=c(2,1))
hist(iris$Sepal.Length)
hist(iris$Sepal.Length)

par(mfrow=c(2,2))
hist(iris$Sepal.Length)
hist(iris$Sepal.Length)
hist(iris$Sepal.Length)
hist(iris$Sepal.Length)

hist(iris$Sepal.Length)
plot(sepal_density)
stripchart(iris$Sepal.Length)
barplot(sepal_table)

par(mfrow=c(1,1))
```

Annotations

```
plot(iris$Petal.Length, iris$Petal.Width)

text(1,2.5, labels="Obscured text")
text(1,2, labels="Text adjusted to 0", adj=0)
text(1,1.5, labels="Text adjusted to 0.5", adj=0.5)
text(1,1, labels="Text adjusted to 1", adj=1)

text(5,0.5, labels="Text with different x position")

text(4,2.5, labels="Obscured text")
text(4,2, labels="Text adjusted to 0", adj=0)
text(4,1.5, labels="Text adjusted to 0.5", adj=0.5)
text(4,1, labels="Text adjusted to 1", adj=1)
```



SCOTTY, I NEED

MORE POWER!

ggplot2

Pros:

- Beautiful (some say)
- Extensible
- Flexible

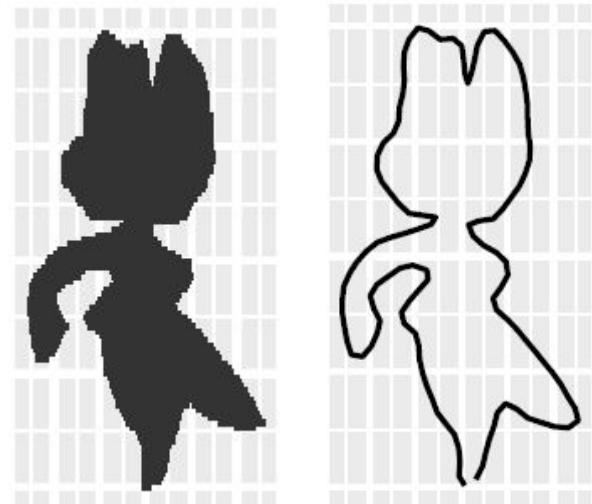
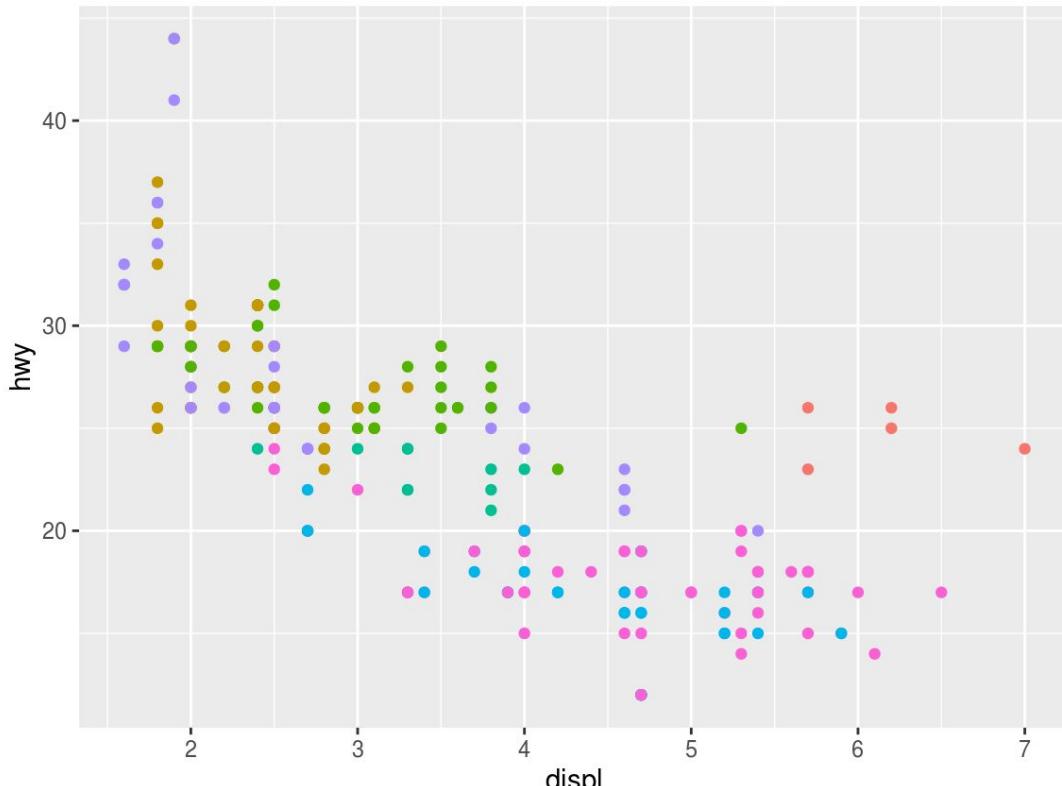
Cons:

- Less convenient
- Different logic can lead to unexpected behaviour

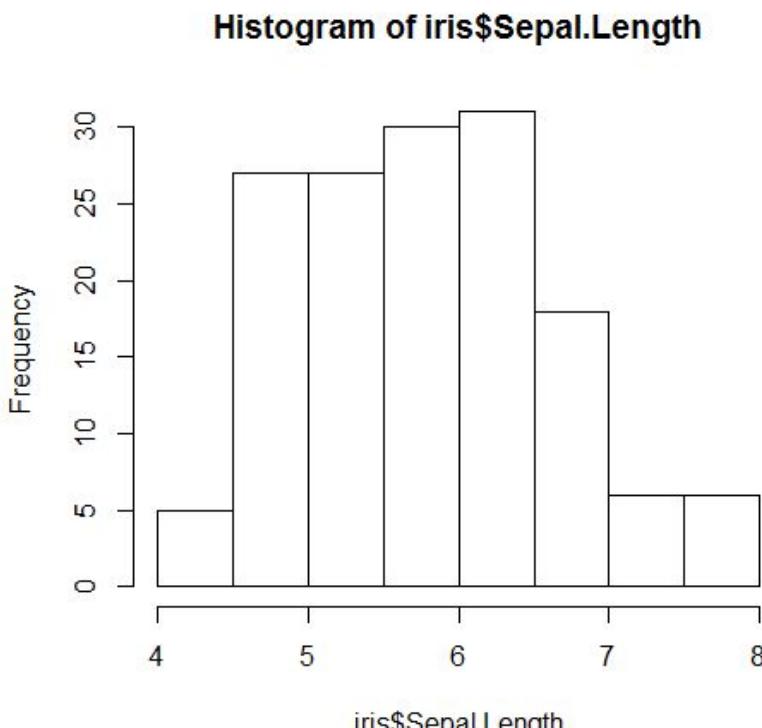
Recommended usage scenario:

Complex or final data visualisation.

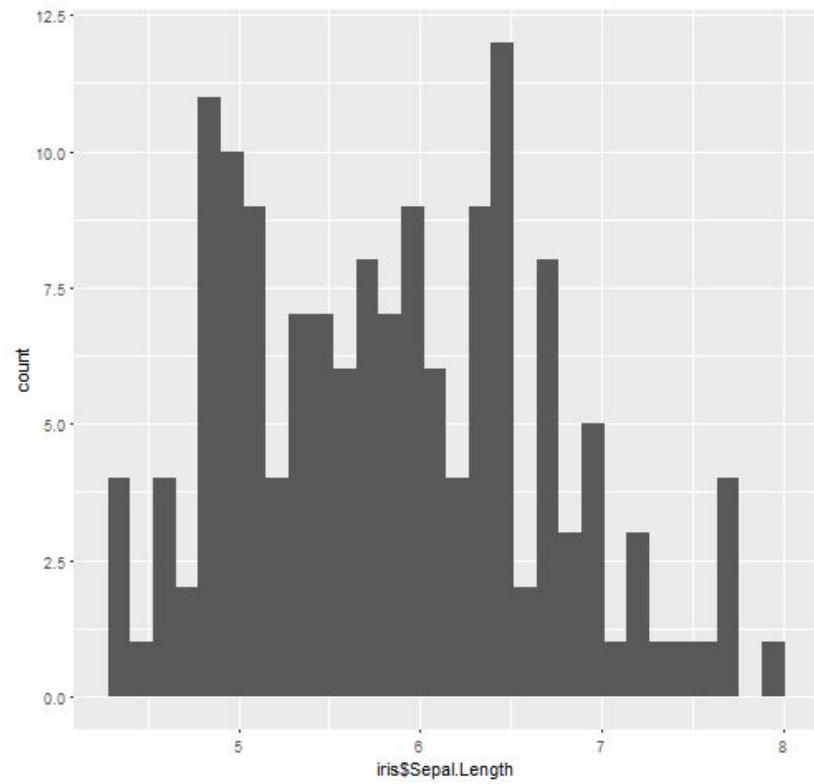
Also fun.



(can be) as easy to use as base r



```
hist(iris$Sepal.Length)
```



```
qplot(iris$Sepal.Length)
```

Exercise 02

Quick recreations of base R plots in ggplots

Single variable plots

- Histograms
- Dot plots
- Strip charts

Multiple variable plots

- Scatterplots
- Paired scatterplots
- Boxplots



Output of other functions or operations

- Density plots
- Linear model diagnostics
- Bar plots
- Stacked bar plots
- Pie charts

Multiple plots

Annotations

Single variable plots

Histogram

```
ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram()
```

Dot chart

```
ggplot(data=iris, aes(x=Sepal.Length)) + geom_dotplot()
```

Strip chart

- No base ggplot version of a strip chart

Multiple variable plots

Scatterplot

```
ggplot(data=iris, aes(x=Petal.Length, y=Petal.Width)) + geom_point()
```

Adding lines to a scatterplot

```
ggplot(data=iris, aes(x=Petal.Length, y=Petal.Width)) + geom_line()
```

Pairs plot

```
library(GGally)
ggpairs(iris[,c("Sepal.Length", "Petal.Width", "Sepal.Width")])
```

Box plot

```
ggplot(data=iris, aes(x=Species, y=Sepal.Length)) + geom_boxplot()
```

Output of other functions or operations

Density plot (cousin of the histogram)

```
ggplot(data=iris, aes(x=Sepal.Length)) + geom_density()
```

Linear model diagnostics

- No base ggplot version of linear model diagnostics

Linear model output on points

```
ggplot(data=iris, aes(x=Sepal.Length, y=Petal.Width)) +  
       geom_point() + geom_smooth(method="lm")
```

Bar plots

```
ggplot(data=iris, aes(x=Sepal.Length)) + geom_bar()
```

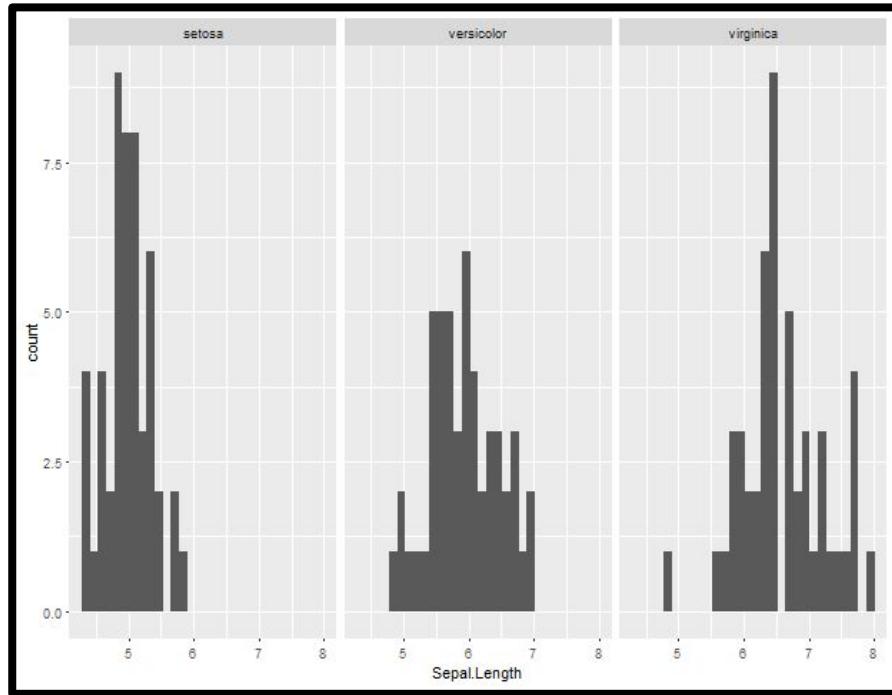
Stacked bar plots

```
ggplot(data=iris, aes(x=Sepal.Length)) +  
       geom_bar(aes(fill=Species))
```

Pie chart

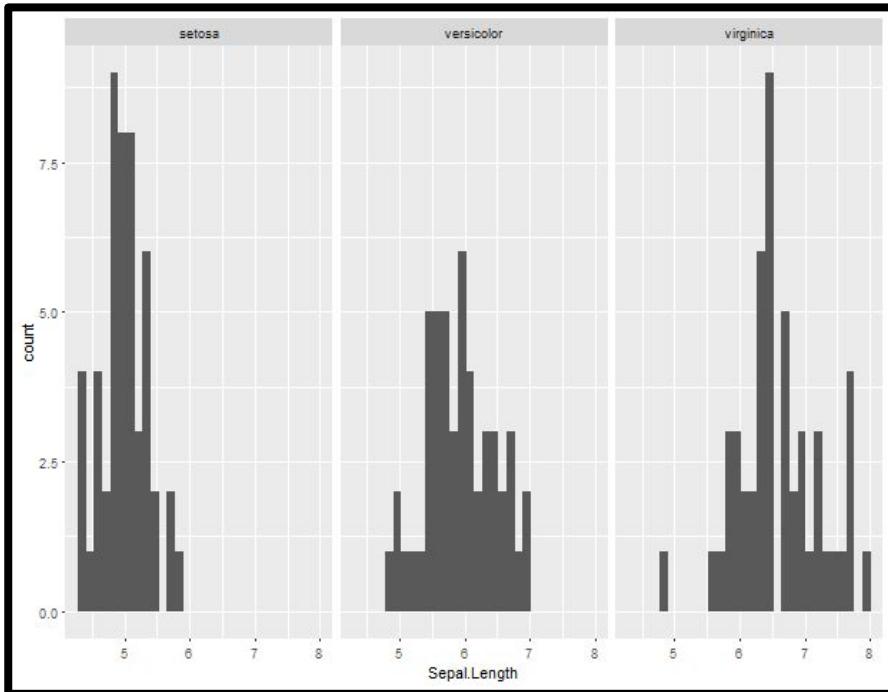
```
ggplot(iris, aes(x="", y=Petal.Width, fill=Species)) +  
       geom_bar(width = 1, stat = "identity") +  
       coord_polar("y", start=0)
```

Multiple plots

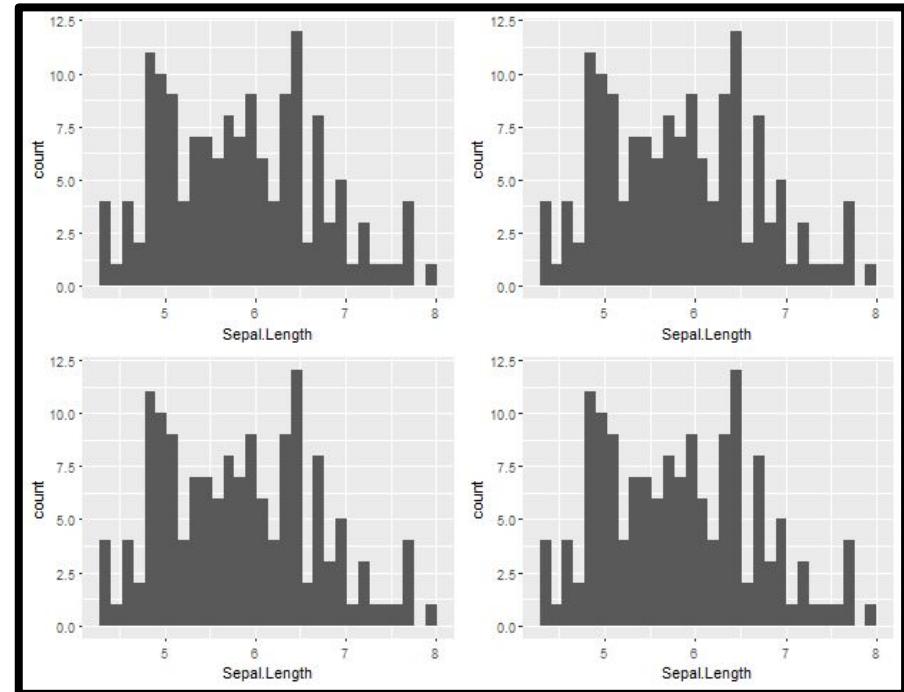


```
ggplot(data=iris, aes(x=Sepal.Length)) +  
  geom_histogram() + facet_wrap(~Species)
```

Multiple plots



```
ggplot(data=iris, aes(x=Sepal.Length)) +  
  geom_histogram() + facet_wrap(~Species)
```



```
library(gridExtra)  
grid.arrange(plot_1, plot_2, plot_3, plot_4)
```

Multiple plots

Facets

```
ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram() +  
  facet_wrap(~Species)
```

Works across almost all plot types

```
ggplot(data=iris, aes(x=Petal.Length, y=Petal.Width)) +  
  geom_point() + facet_wrap(~Species)
```

Still works, even if it is redundant.

```
ggplot(data=iris, aes(x=Species, y=Sepal.Length)) +  
  geom_boxplot() + facet_wrap(~Species)  
ggplot(data=iris, aes(x=Sepal.Length)) +  
  geom_bar(aes(fill=Species)) + facet_wrap(~Species)
```

Multiple plots

gridExtra

```
library(gridExtra)

grid.arrange(ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram())
```

This gives extra flexibility, e.g. specifying **rows** and **columns**

```
grid.arrange(nrow=1,
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram())
```

```
grid.arrange(ncol=1,
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram())
```

Multiple plots

gridExtra

Nesting can be done conveniently via **Grobs**:

```
grid.arrange(ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             arrangeGrob(nrow=1,
                         ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
                         ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
                         ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram()))
```

The major **title** is easily added

```
grid.arrange(top ="This is a title",
             ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
             arrangeGrob(nrow=1,
                         ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
                         ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
                         ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram())
             )
```

Both plots and grobs can be saved as **objects**

```
plot_object<-ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram()
grob_object<- arrangeGrob(nrow=1,
                           ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
                           ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram(),
                           ggplot(data=iris, aes(x=Sepal.Length)) + geom_histogram())

grid.arrange(plot_object,grob_object)
```

Multiple plots

gridExtra

Used in combination with functions and loops,
this is EXTREMELY convenient for iterative plotting.

```
iterative_plot<-list()

for(i in 1:10){iterative_plot[[i]]<-ggplot(data=iris[sample(nrow(iris), 15), ],
                                              aes(x=Petal.Length, y=Petal.Width)) + geom_point()}

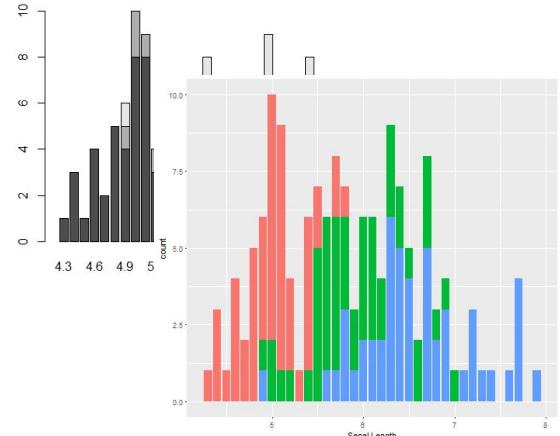
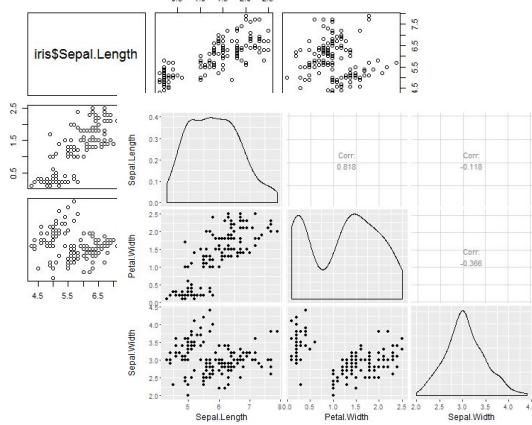
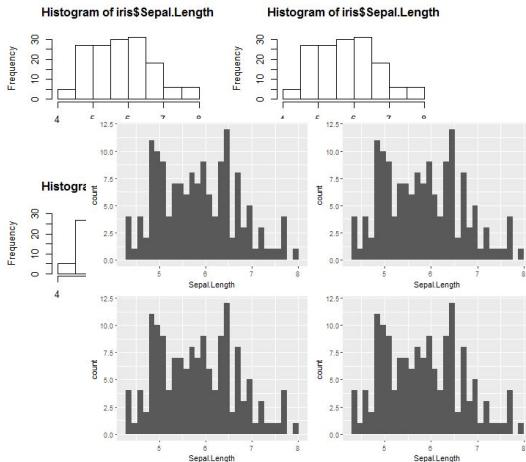
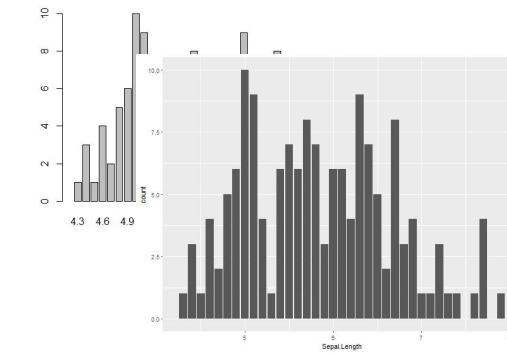
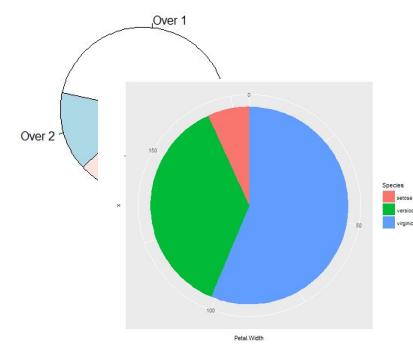
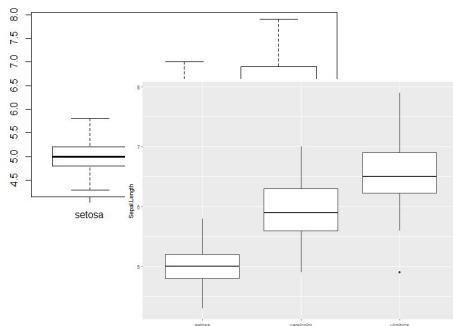
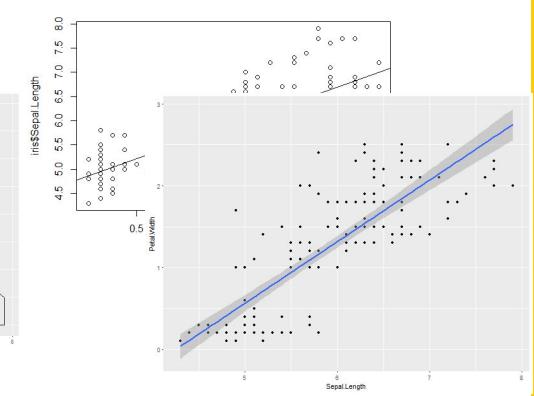
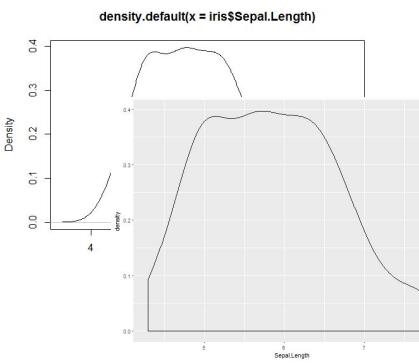
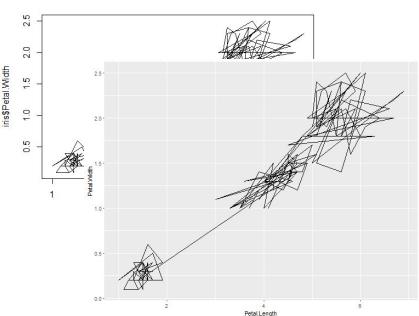
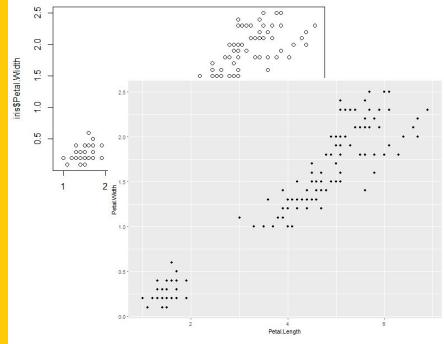
do.call(grid.arrange, iterative_plot)

do.call(grid.arrange, iterative_plot[1:4])

do.call(grid.arrange, c(iterative_plot[1:4], top="Title"))
```

Annotations

```
ggplot(data=iris, aes(x=Petal.Length, y=Petal.Width)) +  
  geom_point() + geom_text(aes(x=5, y=0.5, label="Text"))
```





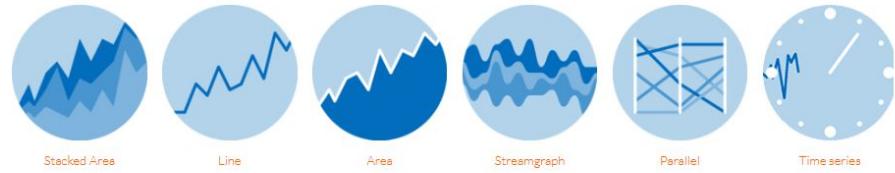
PLEASE █ R

I WANT SOME MORE

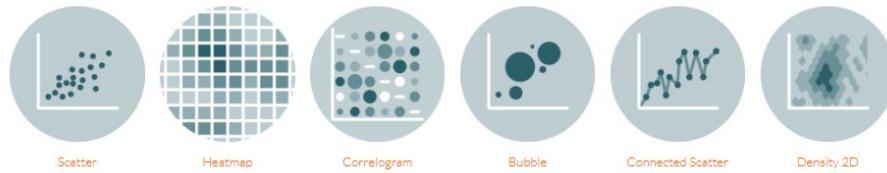
Distribution



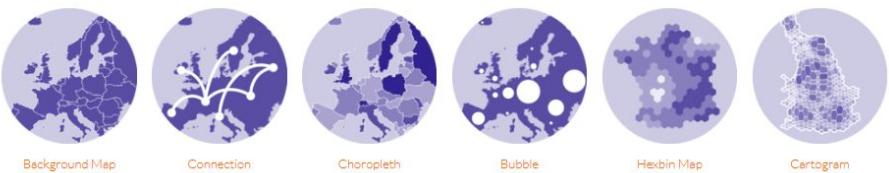
Evolution



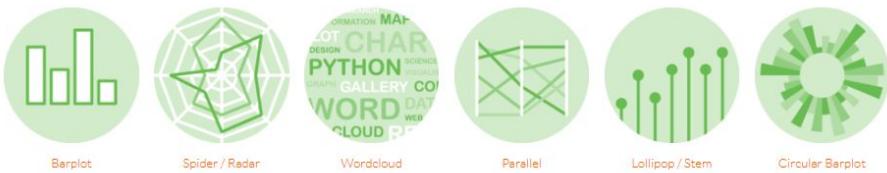
Correlation



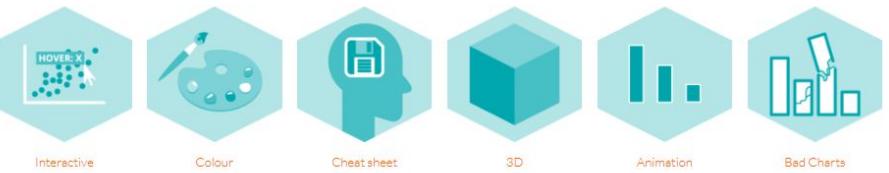
Maps



Rankings



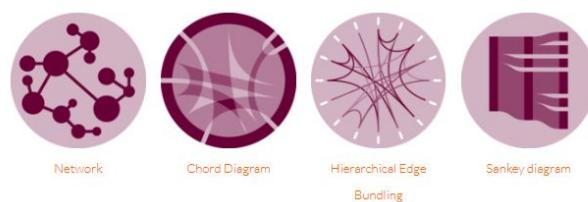
Other



Part of a whole

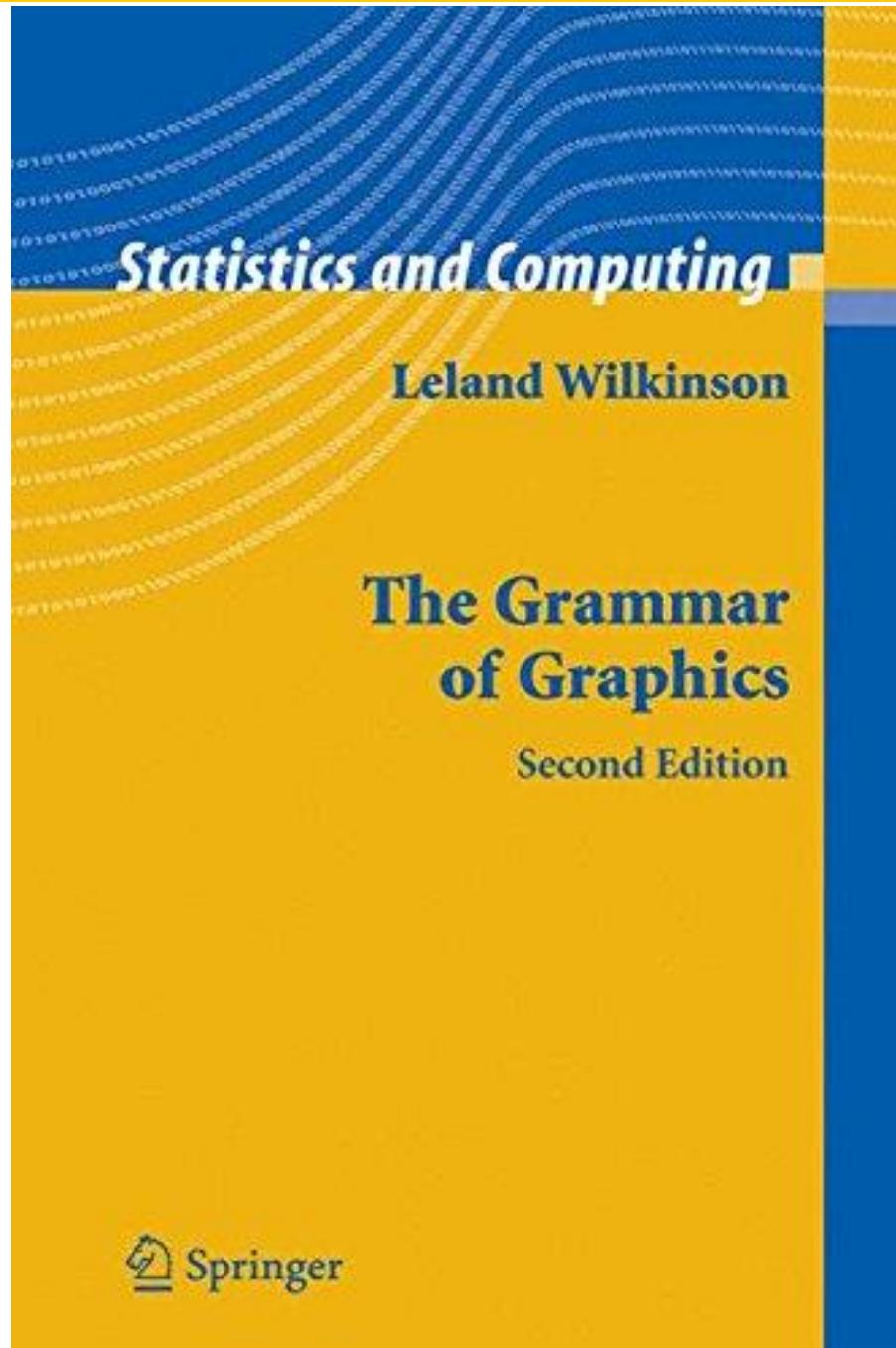


Flow



Examples from <https://www.r-graph-gallery.com/>

Wilkinson, L.
(2006). *The
grammar of
graphics*.
Springer
Science &
Business
Media.



Wilkinson, L. (2006). *The grammar of graphics*. Springer Science & Business Media.

The image shows the front cover of the book 'The Grammar of Graphics' by Leland Wilkinson, Second Edition. The cover has a blue and yellow abstract background with wavy patterns. The title 'Statistics and Computing' is at the top in white, followed by 'Leland Wilkinson' in blue. Below that is the main title 'The Grammar of Graphics' in large blue letters, with 'Second Edition' in smaller blue letters below it. A small circular logo with three overlapping circles is on the left. At the bottom, there is a thin horizontal bar with the text 'Volume 2, November/December 2010' and '© 2010 John Wiley & Sons, Inc.' and the page number '673'.

Advanced Review

The grammar of graphics

Leland Wilkinson*

The grammar of graphics (GoG) denotes a system with seven classes embedded from a raw dataset to a statistical graphic. Each class contains multiple methods, each of which is a function executed at the step in the data flow corresponding to that class. The classes are orthogonal, in the sense that the product set of all classes (every possible sequence of class methods) defines a space of graphics which is meaningful at every point. The meaning of a statistical graphic is thus determined by the mapping produced by the function chain linking method and graphic. © 2010 John Wiley & Sons, Inc. *WIREs Comp Stat* 2010 2: 673–677 DOI: 10.1002/wics.118

Keywords: visualization; statistical graphics

INTRODUCTION

The grammar of graphics (GoG) denotes a system with seven orthogonal classes.¹ The term *orthogonal* means that each class contains one or more methods (functions) as elements, and all tuples in the seven-fold product of these sets of functions produce meaningful graphs. A consequence of this orthogonality is a high degree of expressiveness: we can produce a huge variety of graphical forms or chart types in such a system. In fact, it is claimed that virtually the entire corpus of known charts can be generated by this relatively parsimonious system, and perhaps a great number of meaningful but undiscovered chart types as well.

A second claim of GoG is that this system describes the *meaning* of what we do when we construct statistical graphics. It is more than a taxonomy. It is a computational system based on the underlying mathematics of representing statistical functions of data.

THE GoG DATA FLOW

Figure 1 shows a *data flow* diagram that contains the seven GoG classes. This data flow is a chain that describes the sequence of mappings needed to produce a statistical graphic from a set of data. The first class (*Variables*) maps data to an object called a *varset* (a set of variables). The next two classes (*Algebra*, *Scales*) are transformations on varsets. The next class (*Statistics*) takes a varset and creates a statistical graph

(a statistical summary). The next class (*Geometry*) maps a statistical graph to a geometric graph. The next (*Coordinates*) embeds a graph in a coordinate space. And the last class (*Aesthetics*) maps a graph to a visible or perceivable display called a graphic.

The data flow architecture implies that the subtasks needed to produce a graphic from data must be done in this specified order. Changes to this ordering can produce meaningless graphics. For example, if we compute certain statistics on variables (e.g., sums) before scaling them (e.g., log scales), we can produce statistically meaningless results because the log of a sum is not the sum of the logs.

The data flow in Figure 1 has many paths through it because of the multiple methods (functions) in each stage. We can choose different algebraic designs (factorial, nested, ...), scales (log, probability, ...), statistical methods (means, medians, modes, smoothers, ...), geometric objects (points, lines, bars, ...), coordinate systems (rectangular, polar, ...), and aesthetics (size, shape, color, ...). These paths reveal the richness of the system.

Variables

We begin with data. We assume that the data that we wish to graph are organized in one or more tables. The columns of each table are called *fields*, with each field containing a set of measurements or attributes. The rows of each table are called *records*, with each record containing the measurements of an object on each field. Usually, a relational database management system (RDBMS) produces such a table from organized queries specified in structured query language (SQL) or some other relational language. Other data sources (object, streaming, ...) are mapped to tables through similar methods.

*Correspondence to: Leland.wilkinson@gmail.com
Department of Computer Science, University of Illinois at Chicago, Chicago, IL, USA
DOI: 10.1002/wics.118

Sequential steps in visualisation...

Data goes in

```
1 Sepal.Length Sepal.Width Petal.Length Petal.Width Species  
2 5.1 3.5 1.4 0.2 setosa  
3 4.9 3.0 1.4 0.2 setosa  
4 4.7 3.2 1.3 0.2 setosa  
5 4.6 3.1 1.5 0.2 setosa  
6 5.0 3.6 1.4 0.2 setosa  
7 5.4 3.9 1.7 0.4 setosa
```

Variables

Algebra

Scales

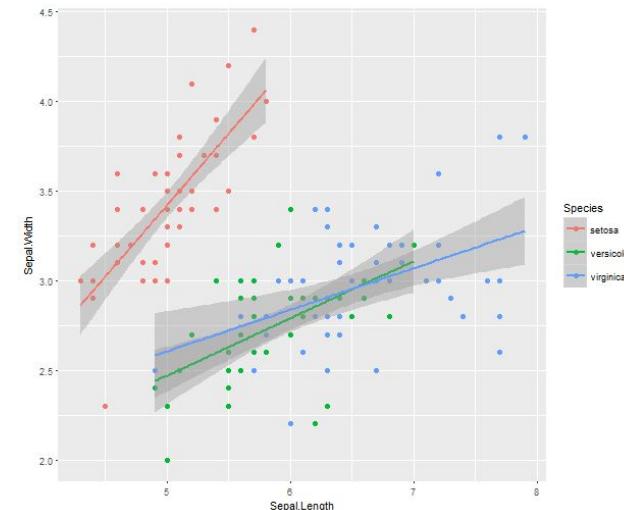
Statistics

Geometry

Coordinates

Aesthetics

Graphics
come out





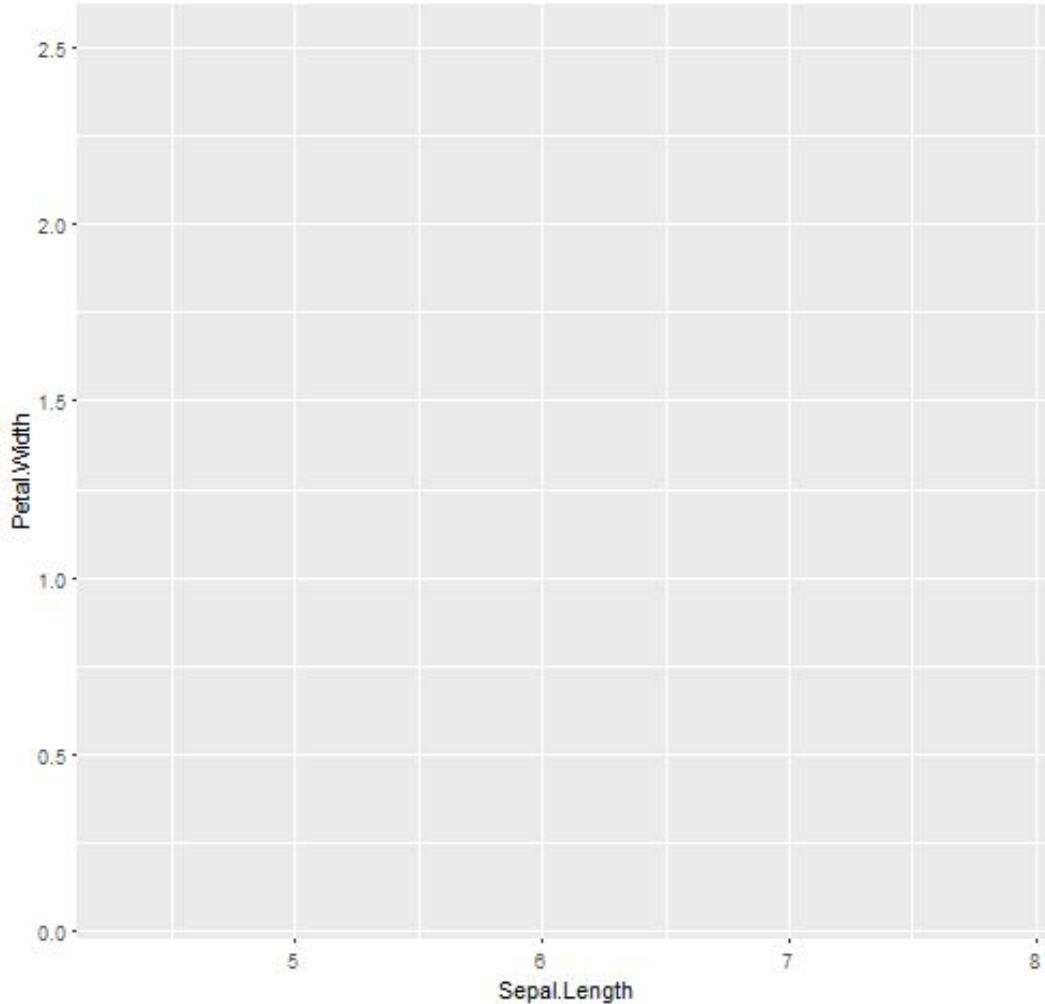
Exercise 03

Follow along, with the `iris` dataset,
looking at `Sepal.Length`
and `Petal.Width`



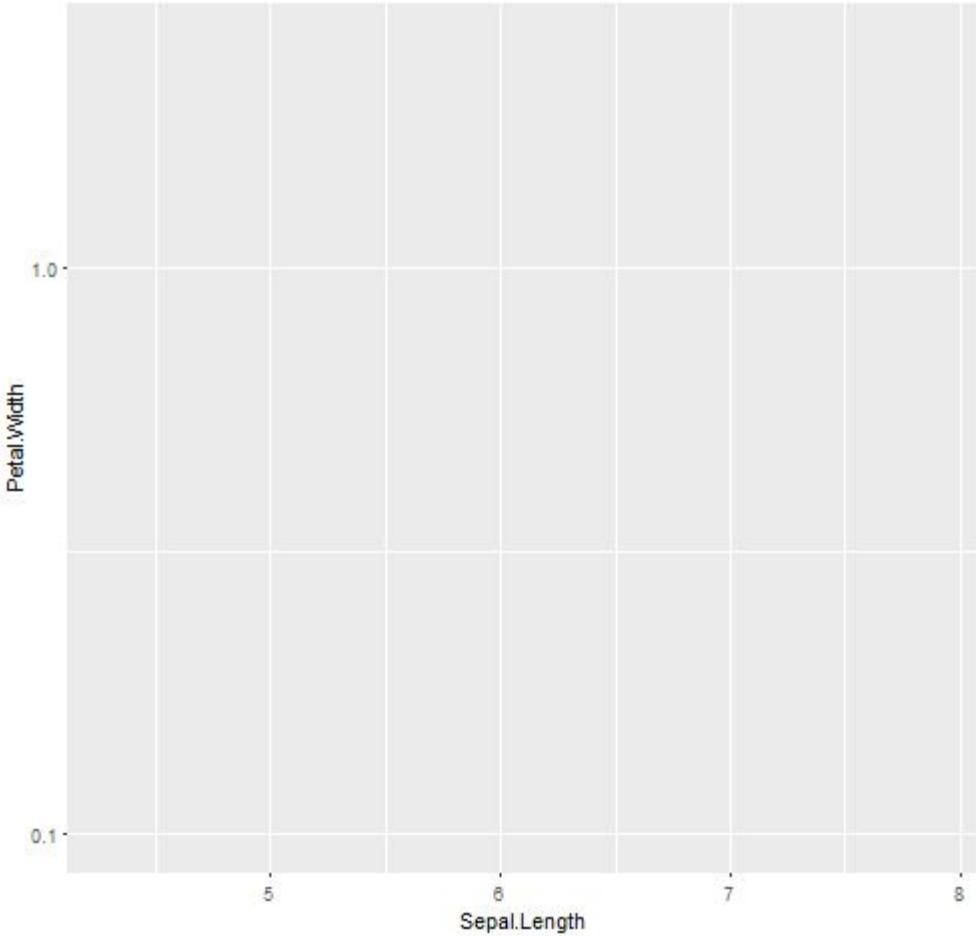


```
ggplot (data=dataset, aes (x=var, y=var2, ...))
```



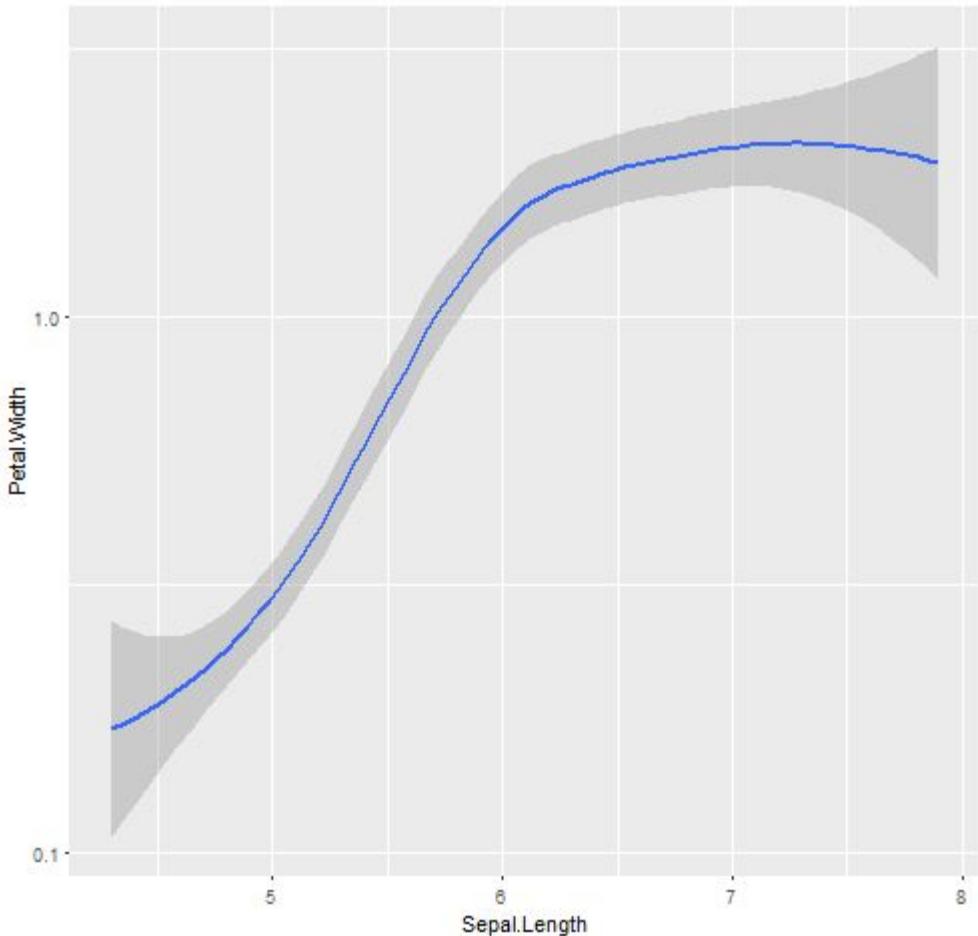


```
ggplot(data=dataset, aes(x=var, y=var2, ...)) +  
  scale_y_log10()
```



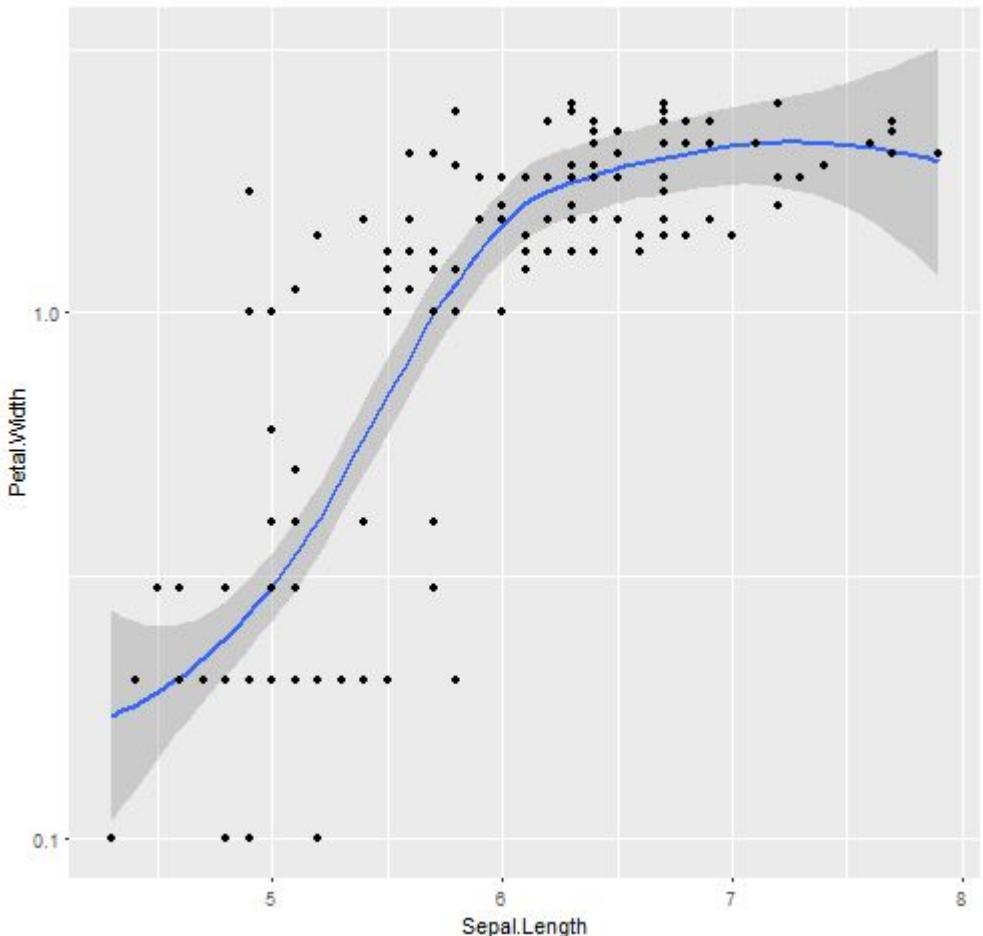


```
ggplot(data=dataset, aes(x=var, y=var2, ...)) +  
  scale_y_log10() +  
  geom_smooth()
```



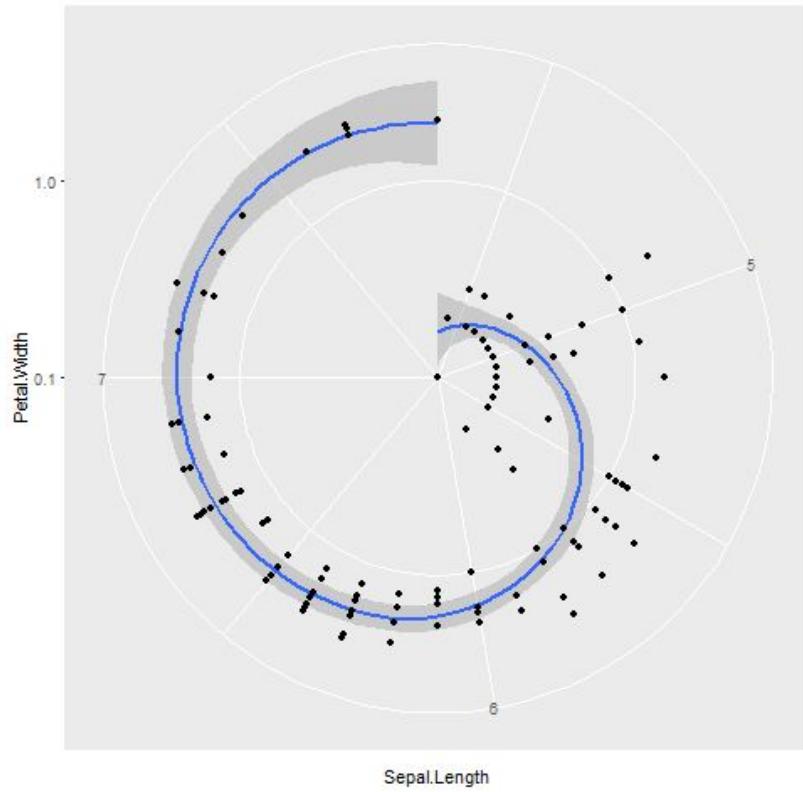


```
ggplot (data=dataset, aes (x=var, y=var2, ...)) +  
  scale_y_log10 () +  
  geom_smooth () +  
  geom_point () +
```





```
ggplot(data=dataset, aes(x=var, y=var2, ...)) +  
  scale_y_log10() +  
  geom_smooth() +  
  geom_point() +  
  coord_polar()
```





```
ggplot(data=dataset, aes(x=var, y=var2, ...)) +  
  scale_y_log10() +  
  geom_smooth(colour="red") +  
  geom_point(colour=yellow) +  
  coord_polar() +  
  theme_bw()
```

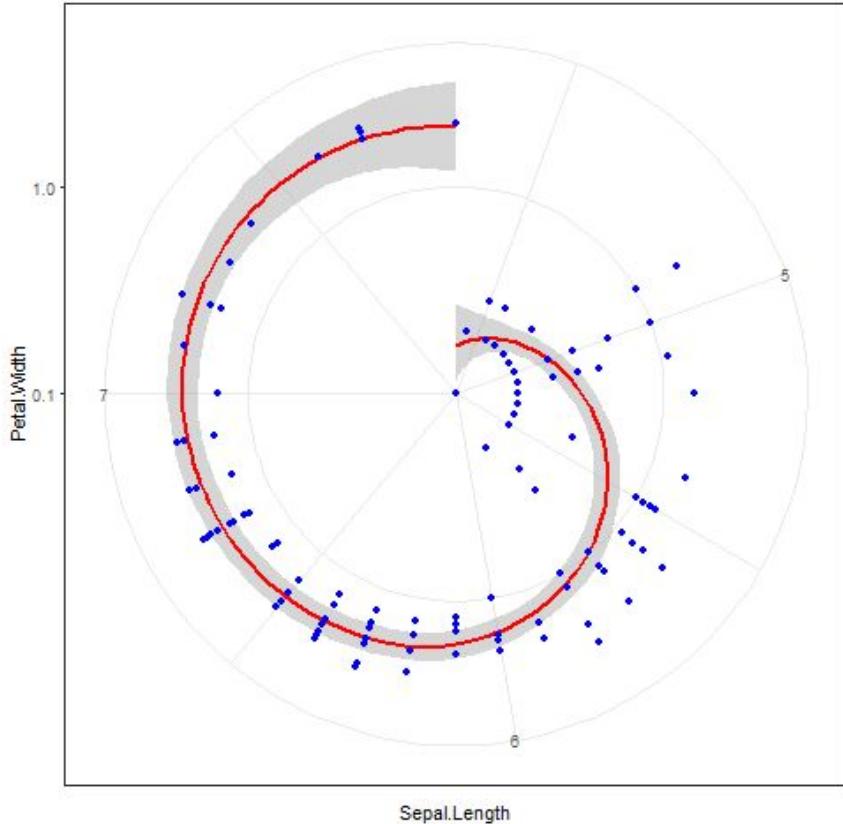




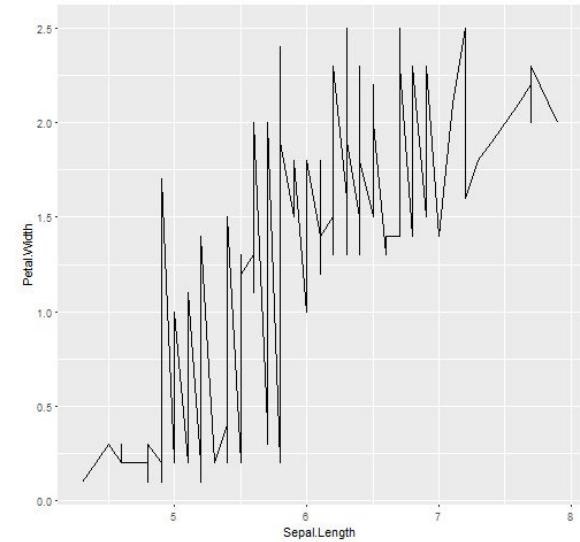
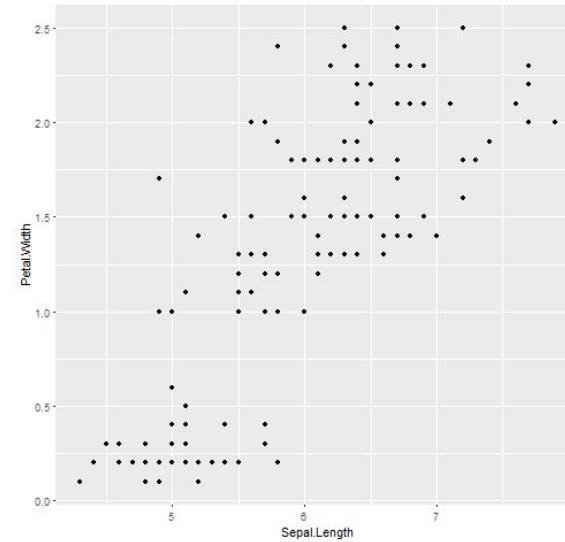
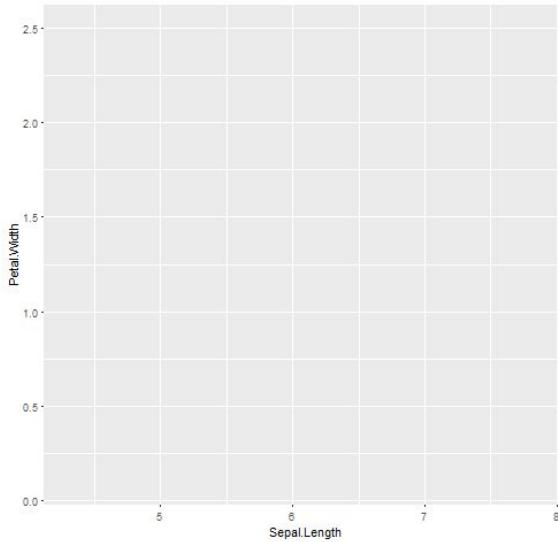
Image source: <https://cityofbyroncity.files.wordpress.com/2012/03/puppy-pics32.jpg>

ggplots are built up in layers

```
base_plot<-ggplot(data=dataset, aes(x=var, y=var2, ...))
```

```
point_plot<-base_plot + geom_point()
```

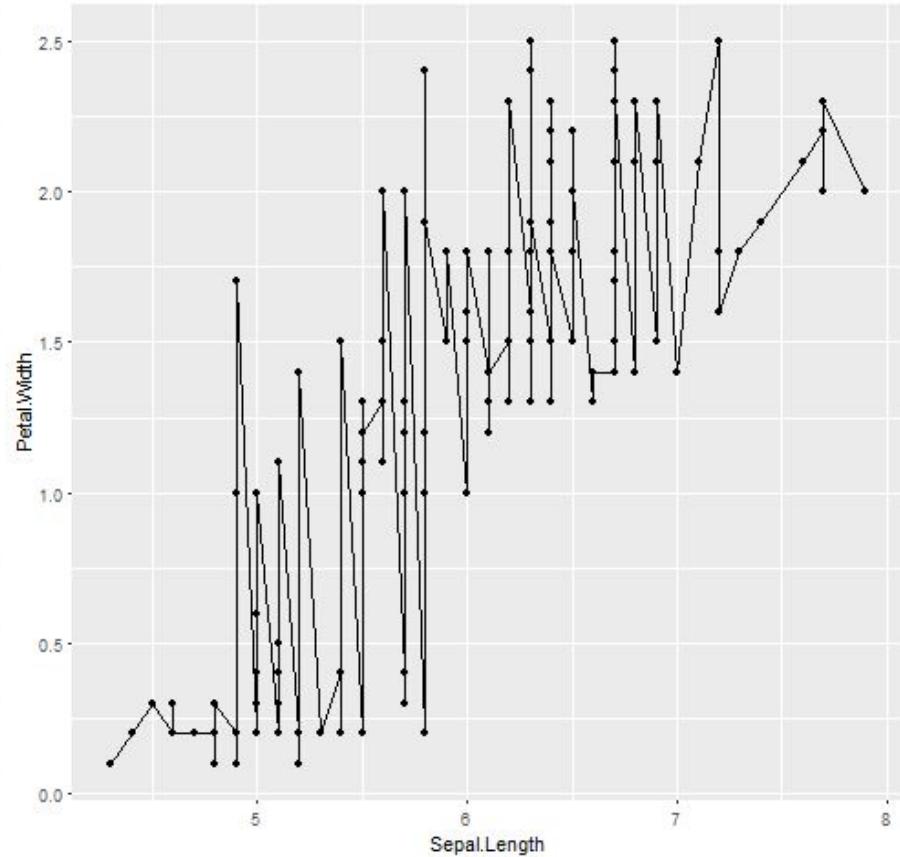
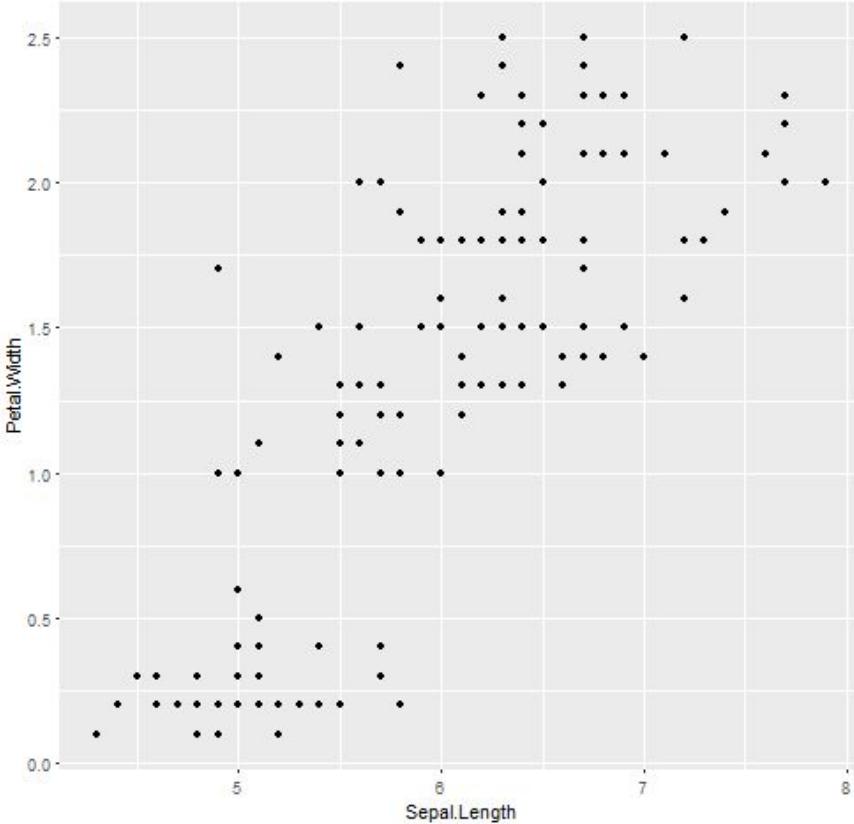
```
line_plot<-base_plot + geom_line()
```



Each step of the grammar can be saved as variables and added together.

```
base_plot + geom_point()
```

```
base_plot + geom_point() + geom_line()
```

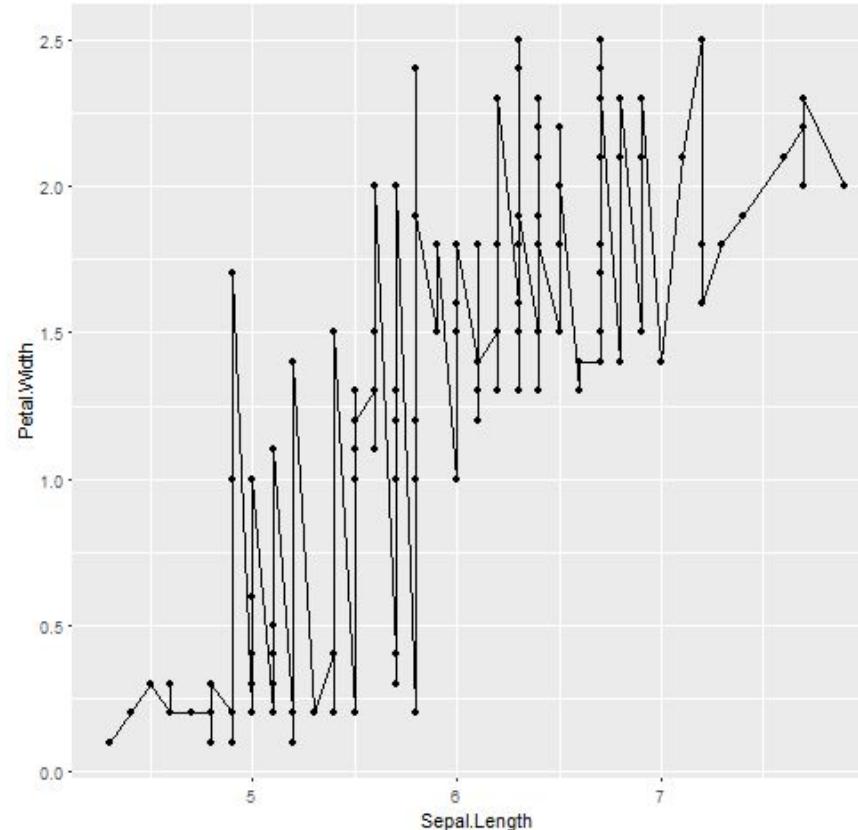


But you can't just add two finished plots.

```
point_plot + line_plot
```

Error: Don't know how to add o to a plot

```
both_plot<-base_plot + geom_point() + geom_line()
```

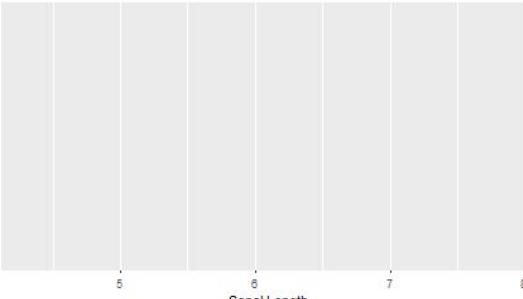


General procedure

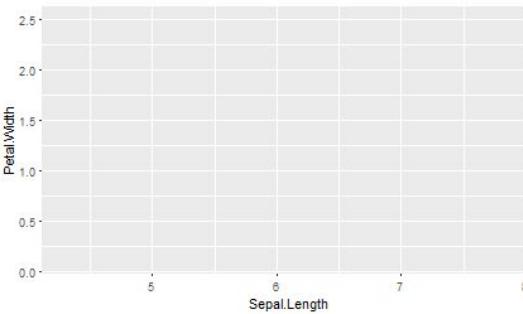
- Make your empty plot with `ggplot()` function
 - This is the point where you specify the data and variables
- Add the actual information with `geom()`
 - You can add multiple layers of data in multiple geoms
- Tweak the axes with `scales()`
- Choose the projection with `coords()`
 - Mostly you skip this step
- Play with the aesthetic within previous steps and use `themes()`

Exercise 04: the plot thickens

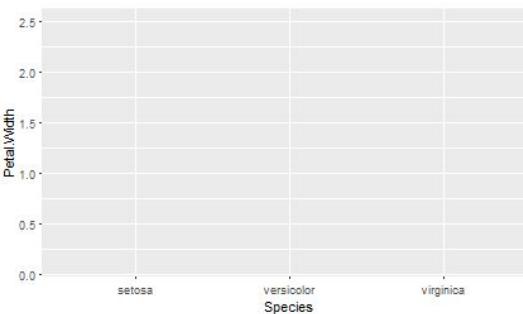
- Go through each step and explore in the Iris data...
- Then build your own plots out of your own data



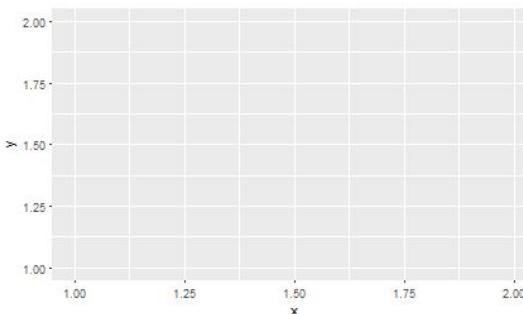
```
base_cont<-ggplot(data=iris,  
aes(x=Sepal.Length))
```



```
base_cont_cont<-ggplot(data=iris,  
aes(x=Sepal.Length, y=Petal.Width))
```



```
base_cont_cat<-ggplot(data=iris,  
aes(x=Species, y=Petal.Width))
```



```
square<-data.frame(x=c(1,2,2,1,1),  
y=c(2,2,1,1,2))  
base_cont_shape<-ggplot(data=square,  
aes(x=x, y=y))
```

Geoms for single continuous variables

`geom_bar()`

`geom_histogram()`

`geom_density()`

`geom_freqpoly()`

Try each one with `base_cont`

Geoms for continuous by continuous variables

```
geom_point()  
geom_jitter()  
geom_line()  
geom_rug()  
geom_step()  
geom_tile()  
geom_bin2d()  
geom_hex()  
geom_quantile()  
geom_smooth()
```

Try each one with `base_cont_cont`

Geoms for continuous by categorical variables

```
geom_boxplot()
```

```
geom_boxplot() + geom_jitter()
```

Try with [base_cont_cat](#)

Further references: <http://sape.inf.usi.ch/quick-reference/ggplot2/geom>

Geoms for shapes

```
geom_path()  
geom_polygon()  
geom_area()
```

Try each one with **base_cont_shape**

Further references: <http://sape.inf.usi.ch/quick-reference/ggplot2/geom>

Geoms for error bars

Ggplots doesn't do this for you, so you need to prepare the data first...

e.g. for error bars reflecting standard deviations by group:

```
summary_grabber<-function(x,group){data.frame(species=group,
                                                 sepal.Length=mean(iris[which(iris$Species==group),x]),
                                                 sd=sd(iris[which(iris$Species==group),x]))}

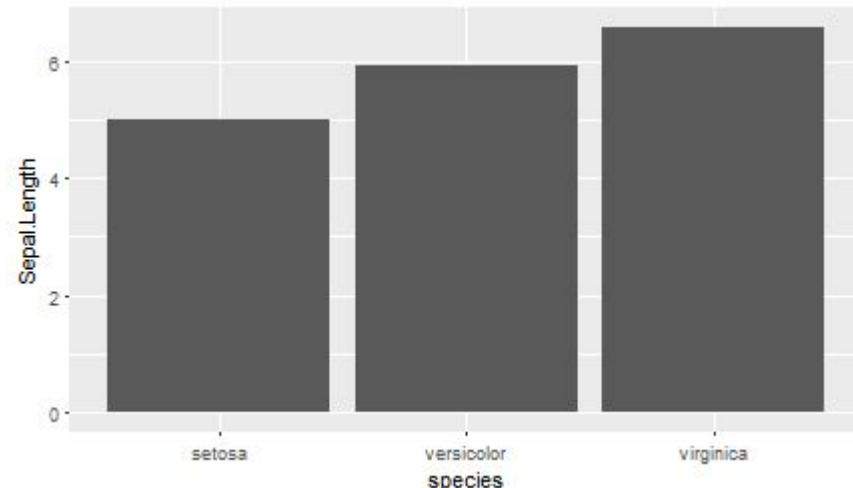
se_dat<-rbind(summary_grabber("Sepal.Length","setosa"),
               summary_grabber("Sepal.Length","versicolor"),
               summary_grabber("Sepal.Length","virginica"))

se_dat$sd_lower<-se_dat$Sepal.Length-se_dat$sd
se_dat$sd_upper<-se_dat$Sepal.Length+se_dat$sd
```

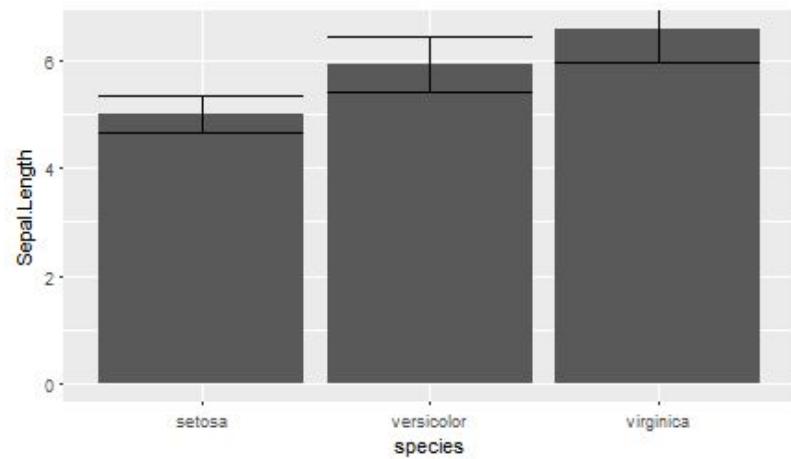
Geoms for error bars

```
> se_dat  
  species Sepal.Length      sd      ymin      ymax sd_lower sd_upper  
1  setosa     5.006 0.3524897 4.653510 5.358490 4.653510 5.358490  
2 versicolor   5.936 0.5161711 5.419829 6.452171 5.419829 6.452171  
3 virginica    6.588 0.6358796 5.952120 7.223880 5.952120 7.223880
```

```
bar_base <- ggplot(data=se_dat,  
                     aes(x=species,  
                          y=Sepal.Length)) +  
  geom_bar(stat="identity")
```



```
bar_base +  
  geom_errorbar(ymin=se_dat$sd_lower,  
                 ymax=se_dat$sd_upper)
```



Geoms for error bars

geom_errorbar()
geom_linerange()
geom_pointrange()
geom_crossbar()

Note: geom_errorbarh() works the same as geom_errorbar() but for the x axis

Try each one with [bar_base](#)

Scales and axes

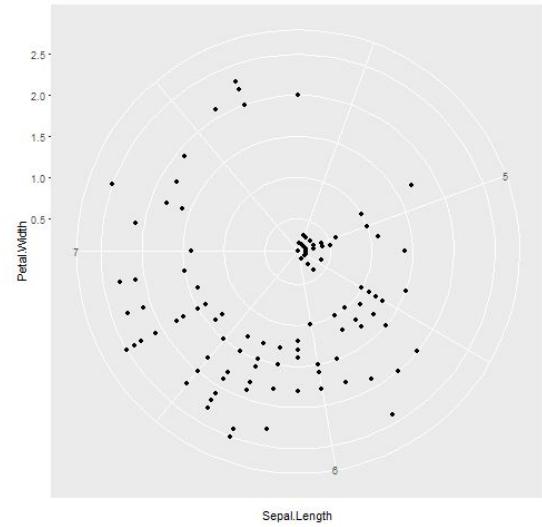
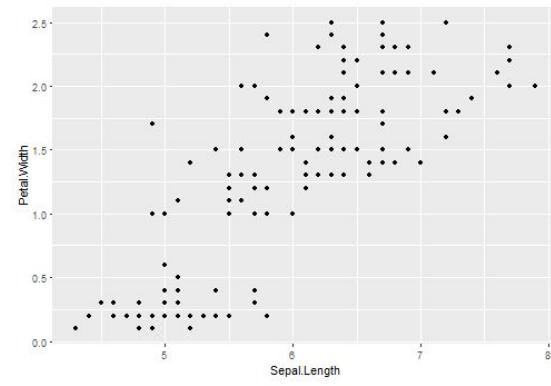
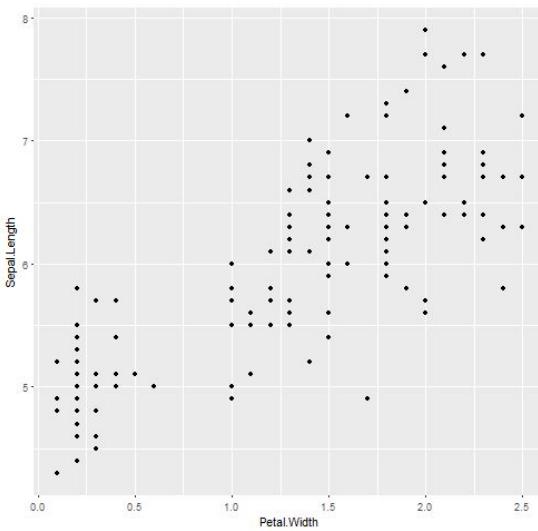
Reverse the axes: `scale_x_reverse()` and `scale_y_reverse()`

- Transform the axes:
 - At the point of plotting:
 - `ggplot(data=dataset, aes(x=log10(var), y=scale(var2))) + geom_point()`
 - After plotting:
 - `scale_x_log10()` and `scale_y_log10()`
 - `scale_x_sqrt()` and `scale_y_sqrt()`
 - `coord_trans(x="log10", y="sqrt")`

Coords

Change the projection

- **coord_flip()** – reverse x and y axes
- **coord_equal()** – fix aspect ratio, like `asp (TRUE)` in base R
- **coord_polar()** – polar coordinates
- Also **coord_map()** for map projections and **coord_cartesian()**, which is just the default.



Try each one!

Adding and changing text

For the plot area, use a geom:

```
geom_text(label="text", x=xpos, y=ypos)
```

Try each of these steps with **bar_base**:

```
geom_text(label="hello")
```

```
geom_text(label="hello", x=1)
```

```
geom_text(label="hello", x=1, y=1)
```

```
geom_text(label="hello", x=1, y=1, colour="white")
```

The title has its own segment;

```
ggttitle("Title text goes here")
```

Changing axis labels is similar:

```
xlab("X axis label")
```

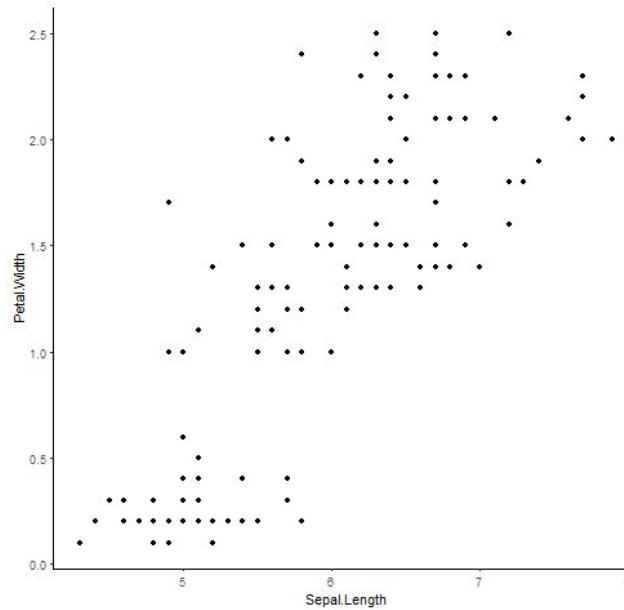
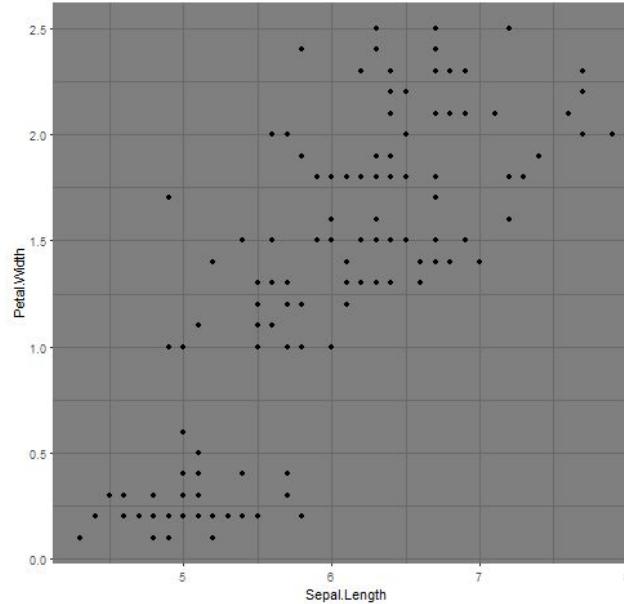
```
ylab("Y axis label")
```

Try to give **bar_base** a unique title and axis labels.

Themes

Many existing themes

- `theme_grey()` (default)
- `theme_bw()`
- `theme_linedraw()`
- `theme_light()`
- `theme_dark()`
- `theme_minimal()`
- `theme_classic()`
- `theme_void()`



Try each one!

Themes

You can make your own, save it as a variable you can add on. Here's just one example to play with:

```
custom_theme <- theme(  
  )
```

Further education:

<https://bookdown.org/rdpeng/RProgDA/building-a-new-theme.html>

Themes

You can make your own, save it as a variable you can add on. Here's just one example to play with:

```
custom_theme <- theme(  
    plot.title = element_text(  
        ) ,  
    panel.grid.major = element_line(  
        ) ,  
    panel.background = element_rect(  
        )  
)
```

Further education:

<https://bookdown.org/rdpeng/RProgDA/building-a-new-theme.html>

Themes

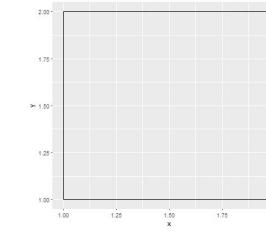
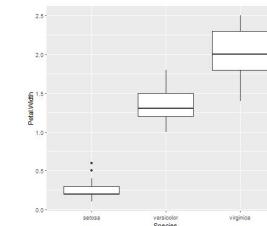
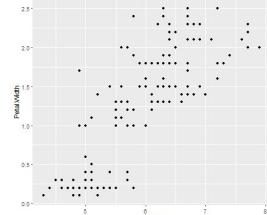
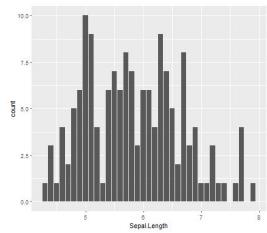
You can make your own, save it as a variable you can add on. Here's just one example to play with:

```
custom_theme <- theme(  
    plot.title = element_text(color = "blue",  
                               face = "bold", size=rel(2)),  
    panel.grid.major = element_line(colour="red",  
                                    linetype = "dotted"),  
    panel.background = element_rect(fill = "yellow")  
)  
  
point_plot + ggtitle("Your own theme") + custom_theme
```

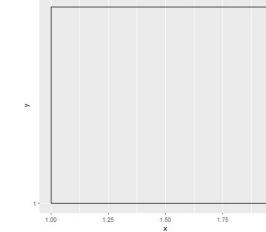
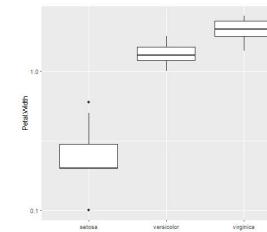
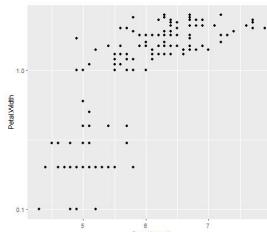
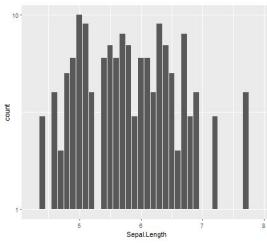
Further education:

<https://bookdown.org/rdpeng/RProgDA/building-a-new-theme.html>

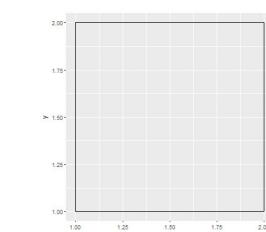
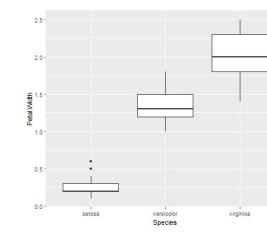
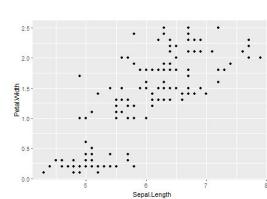
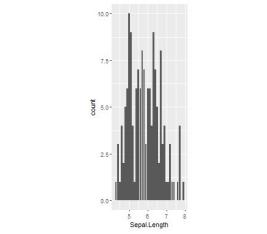
geoms



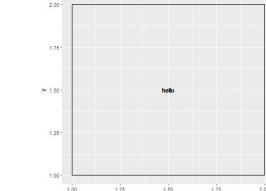
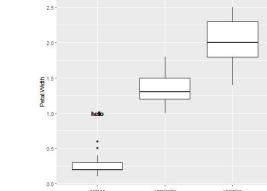
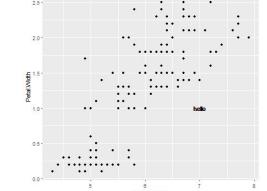
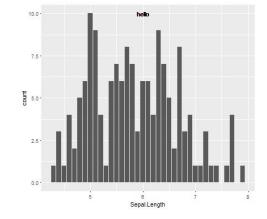
scales and axes scale_y_log()



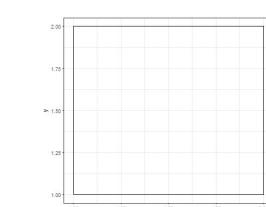
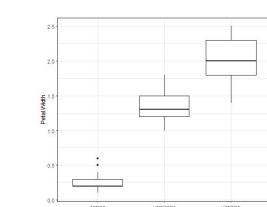
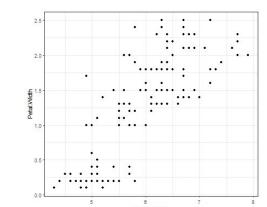
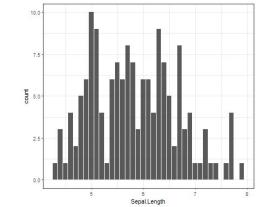
coords coord_equal()



Adding and changing text geom_text()

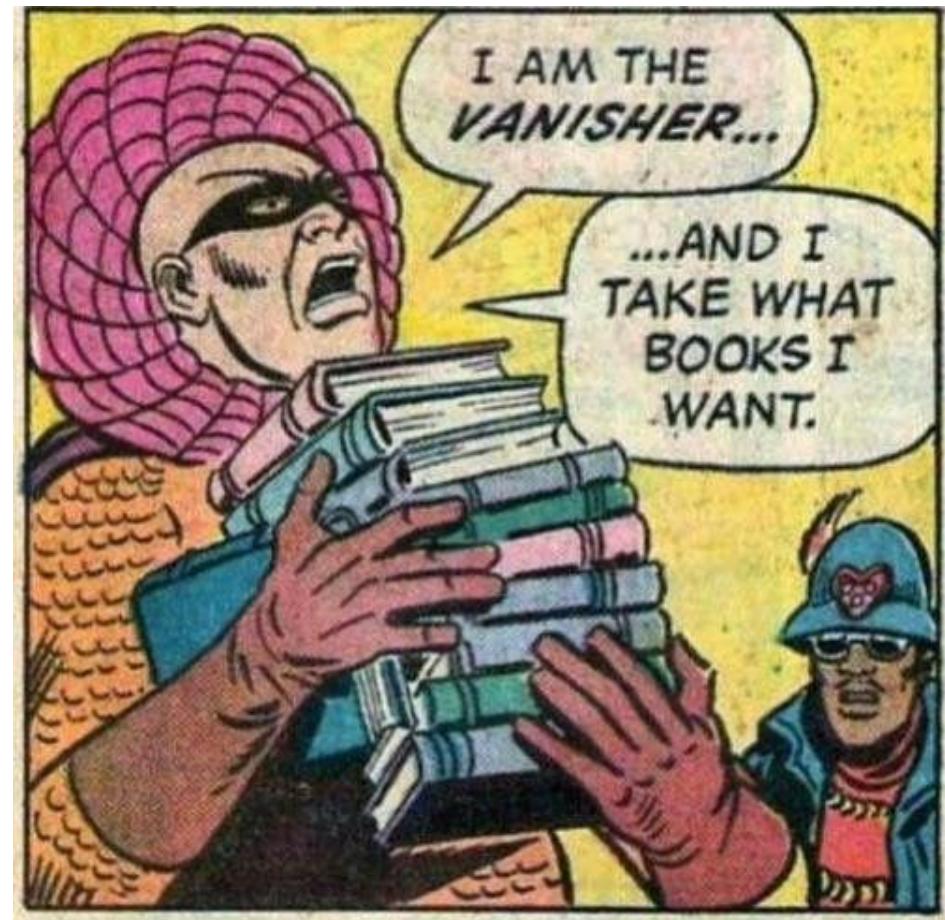


Themes Theme_bw()



Recommended reading

- Tufte, E., & Graves-Morris, P. (2014). *The visual display of quantitative information.*; 1983
- Johnston, S. (2014) R Base graphics: an Idiot's guide
<http://rpubs.com/SusanJohnston/7953>
- Wilkinson, L. (2007) The grammar of graphics. *Journal of statistical software* 17(3). doi: 10.1002/wics.118
- Wickham, H. (2016). *ggplot2: elegant graphics for data analysis*. Springer.
(ANU library has one copy)





Bonus slides: the grammar
of graphics



- Variable = each piece of data we wish to visualise
 - (Varset is a set of one or more variables)



- Variable = each piece of data we wish to visualise
 - (Varset is a set of one or more variables)

• Algebra = combining the varset together.

```
> iris[,c("Sepal.Length", "Sepal.width", "Species")]
```

	Sepal.Length	Sepal.Width	Species	59	6.6	2.9 versicolor	116	6.4	3.2 virginica
1	5.1	3.5	setosa	60	5.2	2.7 versicolor	117	6.5	3.0 virginica
2	4.9	3.0	setosa	61	5.0	2.0 versicolor	118	7.7	3.8 virginica
3	4.7	3.2	setosa	62	5.9	3.0 versicolor	119	7.7	2.6 virginica
4	4.6	3.1	setosa	63	6.0	2.2 versicolor	120	6.0	2.2 virginica
5	5.0	3.6	setosa	64	6.1	2.9 versicolor	121	6.9	3.2 virginica
6	5.4	3.9	setosa	65	5.6	2.9 versicolor	122	5.6	2.8 virginica
7	4.6	3.4	setosa	66	6.7	3.1 versicolor	123	7.7	2.8 virginica
8	5.0	3.4	setosa	67	5.6	3.0 versicolor	124	6.3	2.7 virginica
9	4.4	2.9	setosa	68	5.8	2.7 versicolor	125	6.7	3.3 virginica
10	4.9	3.1	setosa	69	6.2	2.2 versicolor	126	7.2	3.2 virginica
11	5.4	3.7	setosa	70	5.6	2.5 versicolor	127	6.2	2.8 virginica
12	4.8	3.4	setosa	71	5.9	3.2 versicolor	128	6.1	3.0 virginica
13	4.8	3.0	setosa	72	6.1	2.8 versicolor	129	6.4	2.8 virginica
14	4.3	3.0	setosa	73	6.3	2.5 versicolor	130	7.2	3.0 virginica
15	5.8	4.0	setosa	74	6.1	2.8 versicolor	131	7.4	2.8 virginica
16	5.7	4.4	setosa	75	6.4	2.9 versicolor	132	7.9	3.8 virginica
17	5.4	3.9	setosa	76	6.6	3.0 versicolor	133	6.4	2.8 virginica
18	5.1	3.5	setosa	77	6.8	2.8 versicolor	134	6.3	2.8 virginica
19	5.7	3.8	setosa	78	6.7	3.0 versicolor	135	6.1	2.6 virginica
20	5.1	3.8	setosa	79	6.0	2.9 versicolor	136	7.7	3.0 virginica
21	5.4	3.4	setosa	80	5.7	2.6 versicolor	137	6.3	3.4 virginica
22	5.1	3.7	setosa	81	5.5	2.4 versicolor	138	6.4	3.1 virginica
23	4.6	3.6	setosa	82	5.5	2.4 versicolor	139	6.0	3.0 virginica
24	5.1	3.3	setosa	83	5.8	2.7 versicolor	140	6.9	3.1 virginica
25	4.8	3.4	setosa	84	6.0	2.7 versicolor	141	6.7	3.1 virginica
26	5.0	3.0	setosa	85	5.4	3.0 versicolor	142	6.9	3.1 virginica
27	5.0	3.4	setosa	86	6.0	3.4 versicolor	143	5.8	2.7 virginica
28	5.2	3.5	setosa	87	6.7	3.1 versicolor	144	6.8	3.2 virginica
29	5.2	3.4	setosa	88	6.3	2.3 versicolor	145	6.7	3.3 virginica
30	4.7	3.2	setosa	89	5.6	3.0 versicolor	146	6.7	3.0 virginica
31	4.8	3.1	setosa	90	5.5	2.5 versicolor	147	6.3	2.5 virginica
32	5.4	3.4	setosa	91	5.5	2.6 versicolor	148	6.5	3.0 virginica
33	5.2	4.1	setosa	92	6.1	3.0 versicolor	149	6.2	3.4 virginica
34	5.5	4.2	setosa	93	5.8	2.6 versicolor	150	5.9	3.0 virginica
				94	5.0	2.3 versicolor			
				95	5.6	2.7 versicolor			



- Variable = each piece of data we wish to visualise
 - (Varset is a set of one or more variables)
- Algebra = combining the varset together.
- **Scales** = identifying whether it is nominal, ordinal, interval, or ratio (because each one should be plotted differently)

```
> head(iris[,c("Sepal.Length", "Sepal.width", "species")])
```

	Sepal.Length	Sepal.width	Species
1	5.1	3.5	setosa
2	4.9	3.0	setosa
3	4.7	3.2	setosa
4	4.6	3.1	setosa
5	5.0	3.6	setosa
6	5.4	3.9	setosa

Ratio Ratio Ordinal
 (continuous) (continuous) (categorical)

Fun Tip: The order in the factor will be the order in the plot: to change the order in the plot, change the order of the factor AT THIS POINT.



- Variable = each piece of data we wish to visualise
 - (Varset is a set of one or more variables)
- Algebra = combining the varset together.
- Scales = identifying whether it is nominal, ordinal, interval, or ratio (because each one should be plotted differently)
- **Statistics = receives a varset, computes statistical summaries, and outputs another varset (e.g. means, medians, correlations...)**

```

> summary(glm(sepal.Length ~ Sepal.Width + Species, dat=iris))

Call:
glm(formula = sepal.Length ~ Sepal.Width + Species, data = iris)

Deviance Residuals:
    Min      1Q  Median      3Q      Max 
-1.30711 -0.25713 -0.05325  0.19542  1.41253 

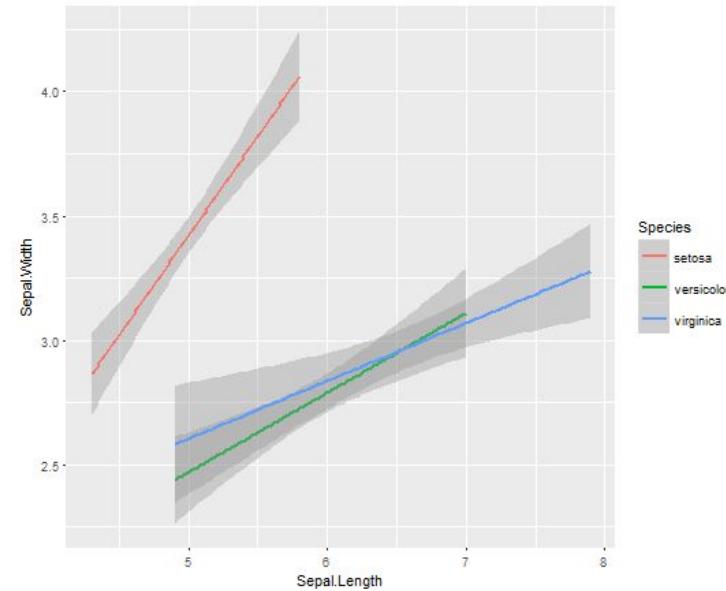
Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept)  2.2514    0.3698   6.089 9.57e-09 ***
Sepal.Width   0.8036    0.1063   7.557 4.19e-12 ***
Speciesversicolor  1.4587    0.1121  13.012 < 2e-16 ***
Speciesvirginica   1.9468    0.1000  19.465 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 0.1918059)

Null deviance: 102.168 on 149 degrees of freedom
Residual deviance: 28.004 on 146 degrees of freedom
AIC: 183.94

Number of Fisher Scoring iterations: 2

```



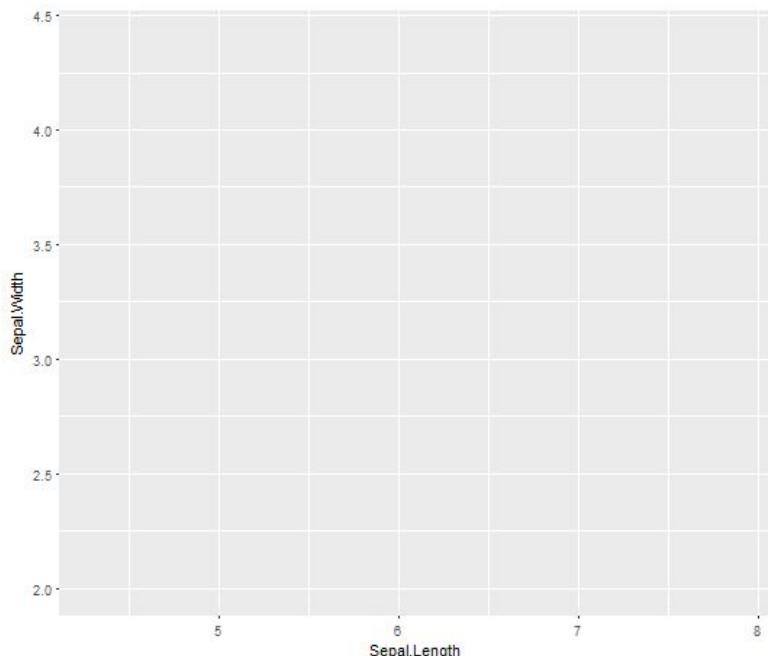
Fun Tip: For most statistics, you can do this manually beforehand, or let ggplot do it for you.



- Variable = each piece of data we wish to visualise
 - (Varset is a set of one or more variables)
- Algebra = combining the varset together.
- Scales = identifying whether it is nominal, ordinal, interval, or ratio (because each one should be plotted differently)
- Statistics = receives a varset, computes statistical summaries, and outputs another varset (e.g. means, medians, correlations...)

• **Geometry** = graphing functions injectively map an m -dimensional bounded region to an n -dimensional real space

AKA the extent of the plot



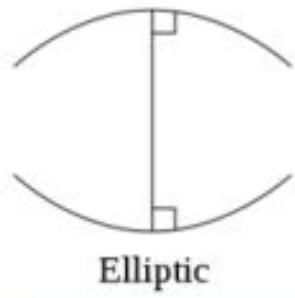
Fun Tip: ggplot plots the data AFTER this. So if you manually change the axes, you change the plot fundamentally (not just plotting everything and shrinking the window down).



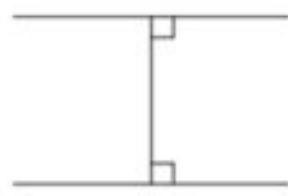
- Variable = each piece of data we wish to visualise
(Varset is a set of one or more variables)

- Algebra = combining the varset together.
- Scales = identifying whether it is nominal, ordinal, interval, or ratio (because each one should be plotted differently)
- Statistics = receives a varset, computes statistical summaries, and outputs another varset (e.g. means, medians, correlations...)
- Geometry = graphing functions injectively map an m -dimensional bounded region to an n -dimensional real space

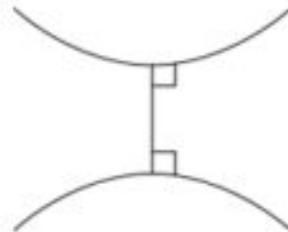
• Coordinates = Cartesian coordinates. 'nuff said.



Elliptic



Euclidean



Hyperbolic

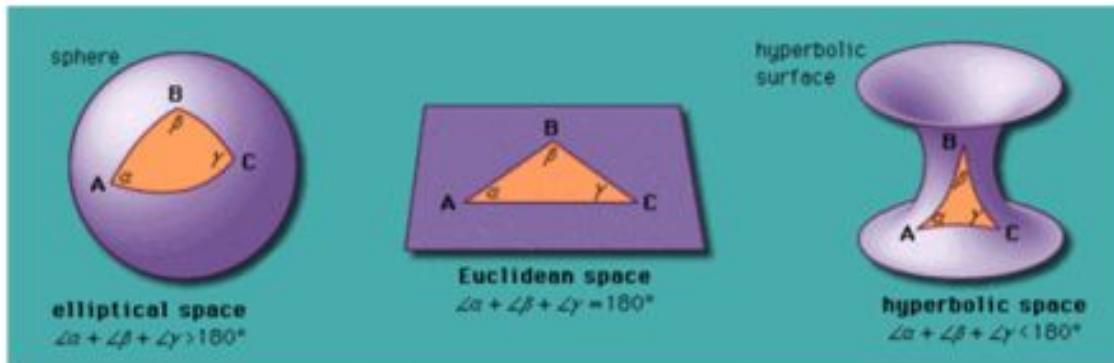
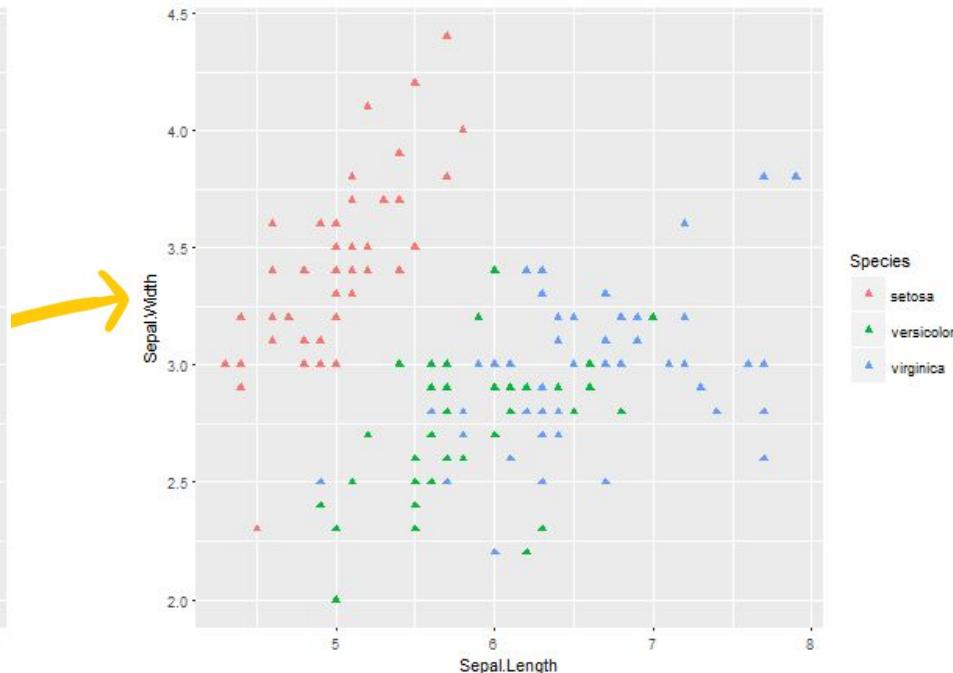
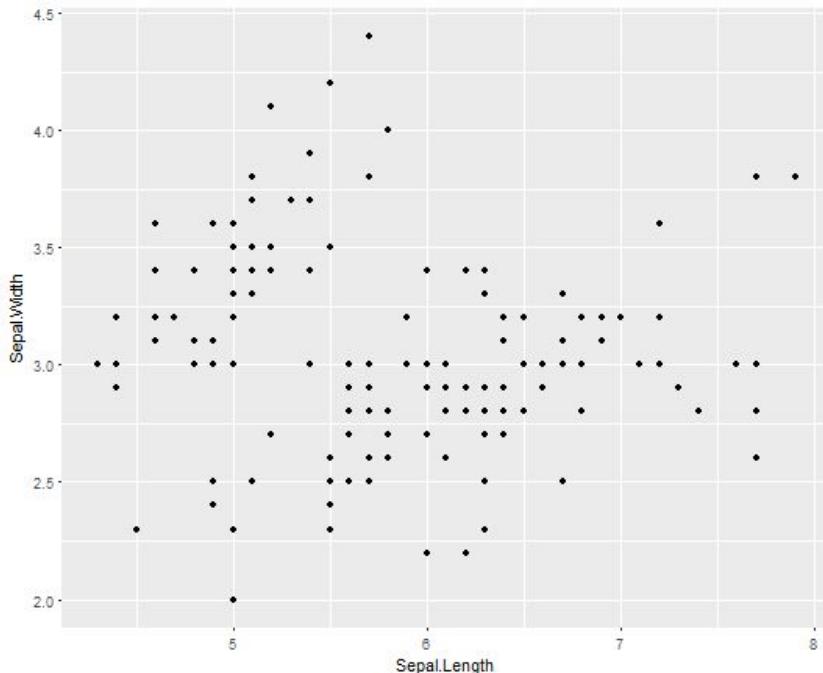


Image source:
<https://naiadseye.files.wordpress.com/2014/10/euclidean-understand-e141449051530.png?w=470&h=289>



- Algebra = combining the varset together.
- Scales = identifying whether it is nominal, ordinal, interval, or ratio (because each one should be plotted differently)
- Statistics = receives a varset, computes statistical summaries, and outputs another varset (e.g. means, medians, correlations...)
- Geometry = graphing functions injectively map an m -dimensional bounded region to an n -dimensional real space
- Coordinates = Cartesian coordinates. 'nuff said.

• Aesthetics = colours, size, shapes,





... But every step is customisable!

```

> iris[,c("Sepal.Length","Sepal.Width","Species")]
   Sepal.Length Sepal.Width Species
1          5.1        3.5  setosa
2          4.9        3.0  setosa
3          4.7        3.2  setosa
4          4.6        3.1  setosa
5          5.0        3.6  setosa
6          5.4        3.9  setosa
7          4.6        3.4  setosa
8          5.0        3.4  setosa
9          4.4        3.0  setosa
10         4.9        3.4  setosa
11         5.4        3.4  setosa
12         4.8        3.0  setosa
13         4.8        3.4  setosa
14         4.3        3.1  setosa
15         5.8        3.0  setosa
16         5.7        3.0  setosa
17         5.4        3.9  setosa
18         5.1        3.3  setosa
  
```

