

Q1: What type of data do we have?

IDA = Initial Data Analysis

- background
- structure
- wrangling
- summaries

main qualities of data suggests population

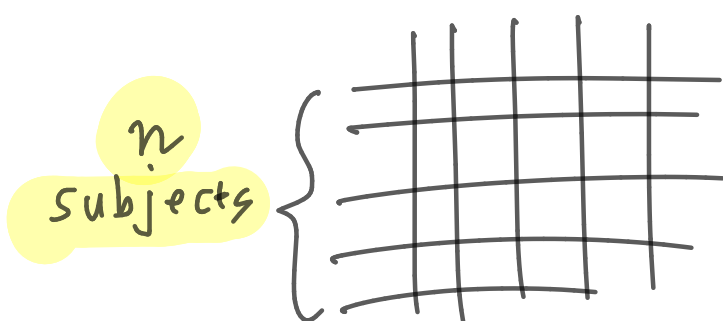


Research Questions



(i) Size (tidy data)

p Variables = attributes of subjects



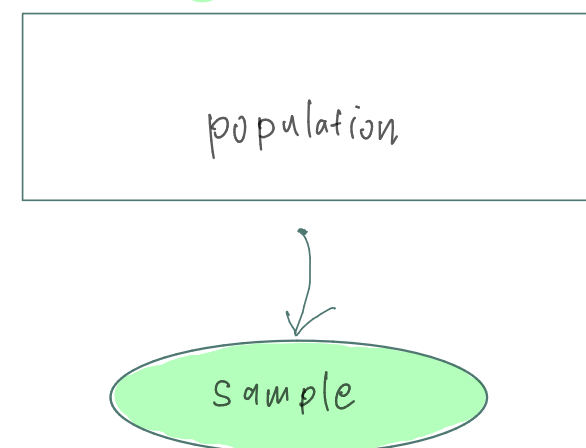
p = dimension

multivariate (2^+)

bivariate (2)

univariate (1)

Module 1: Exploring Data



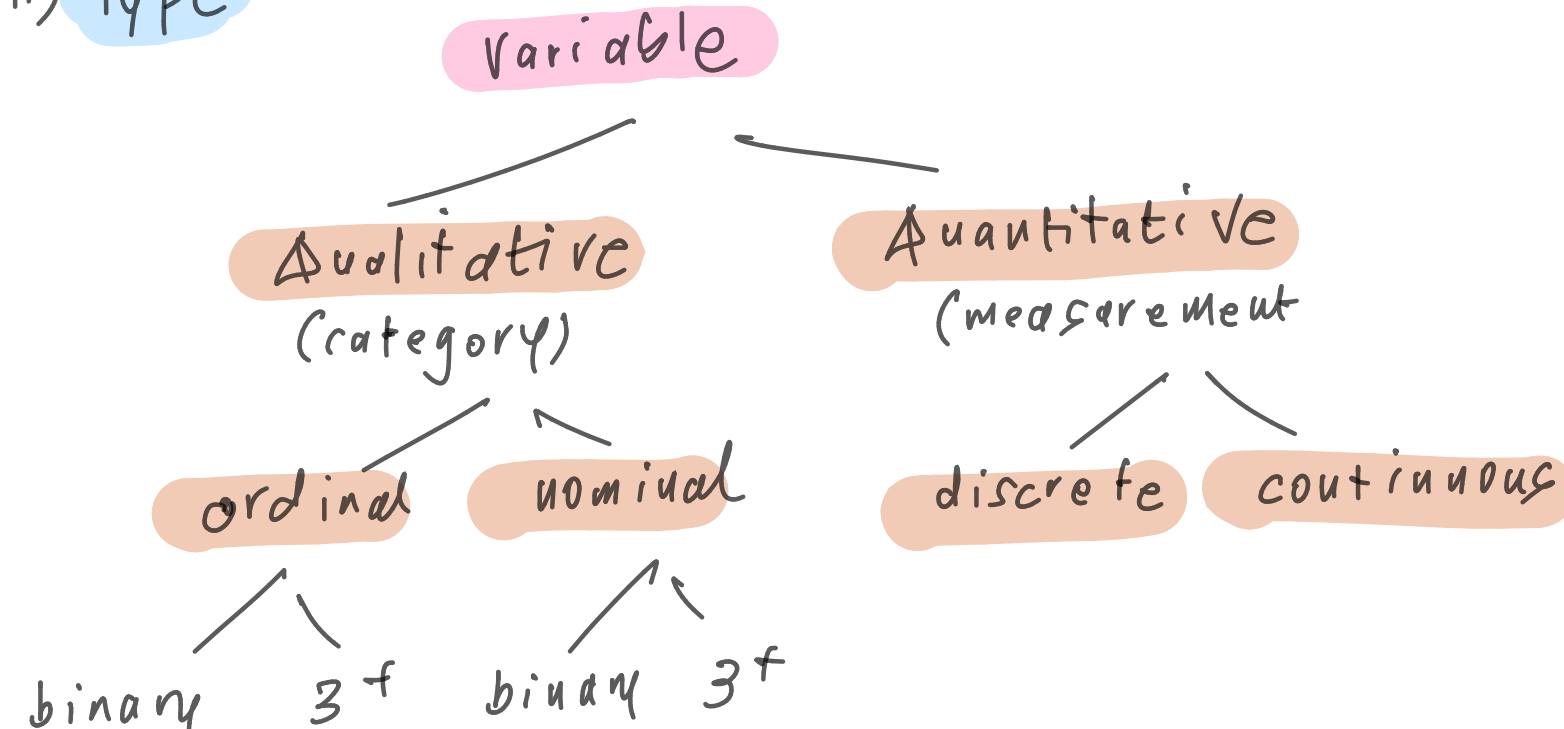
T1: Design of Experiments

T2: Data & Graphical Summaries

T3: Numerical Summaries

LO3: Produce, interpret & compare graphical & numerical summaries using base R & ggplot.

(ii) Type

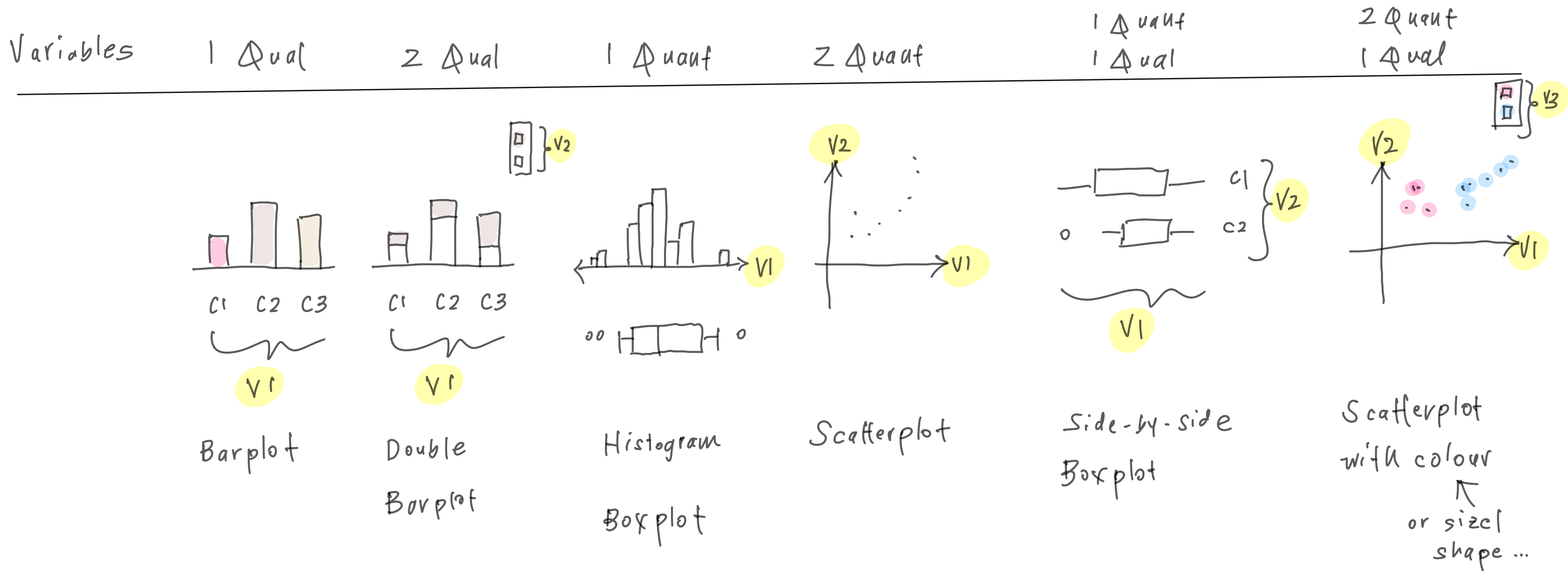


Note:

- ID ≠ Variable DV4/21
- Variables can be Qual or Quant DV4/23

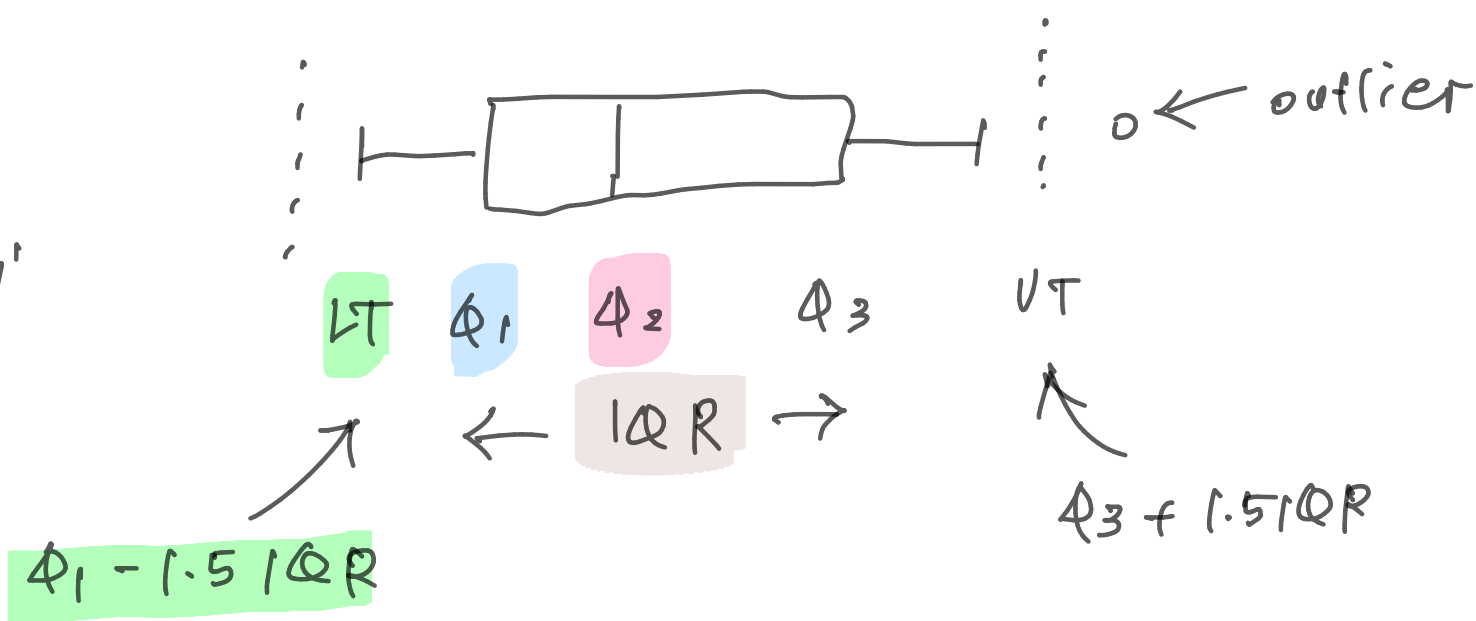
Q2: How can we visualise data?

Data viz ✓ good. = tells an interesting story in a visually appealing way.
 ✗ poor = "chart junk" = visually boring, or distracting
 ✗ bad = misleading story



Note:

- Histograms of Barplots DV5/12-15
- Mistakes with histograms DV5/21-22
- Boxplot is a visual 'numerical summary' DV5/24, DV8



LT = lower threshold

Q1 = 1st Quartile (25% data below)

Q2 = 2nd Quartile (50% data below) = Median

Q3 = 3rd Quartile (75% data below)

IQR = Interquartile Range (50% data between)

How to produce data viz in RStudio?

The **iris** data is already in RStudio!

iris

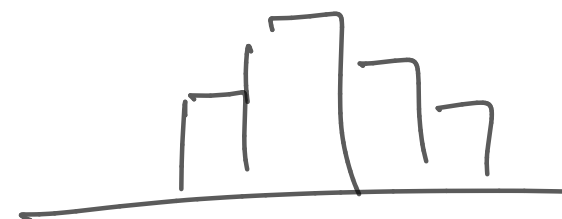
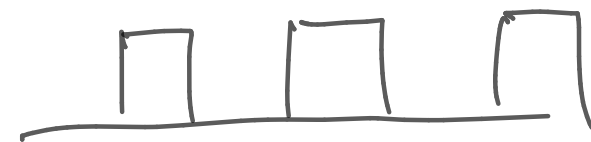
Quant				Qual
Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species

150x5

Base R

```
barplot ( table ( iris $ Species ) )
```

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



ggplot

```
library (tidyverse)
```

```
ggplot ( iris, aes ( x = Species ) )  
  + geom_bar ( )
```

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