

Data analysis methods in weather and climate research

Dr. David B. Stephenson
Department of Meteorology
University of Reading
www.met.rdg.ac.uk/cag



2. Exploratory Data Analysis (EDA)

- Data tabulation
- Summary plots
- Measures of location, scale, and shape
- Rank statistics and empirical quantiles
- Transformation of data values

(c) 2005 D.B.Stephenson@reading.ac.uk

1

2. Data tabulation

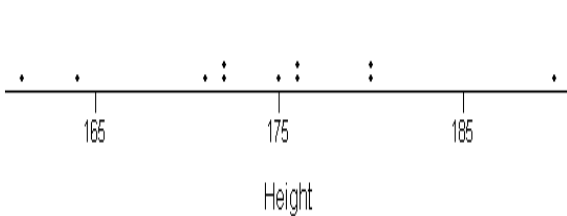
Object	Age (years)	Height (cm)	Weight (kgs)
1	30.9	180	76
2	26.9	164	64
3	33.2	176	87
4	28.5	172	75
5	32.3	176	75
6	37.0	180	86
7	38.3	171	65
8	31.5	172	76
9	32.8	161	75
10	37.7	175	85
11	29.1	190	83

rdgmorph.txt
Meteorologist data

Rows=objects
Columns=variables
Sample size=n=11

2

2. Dotplot (=1-d scatter plot)



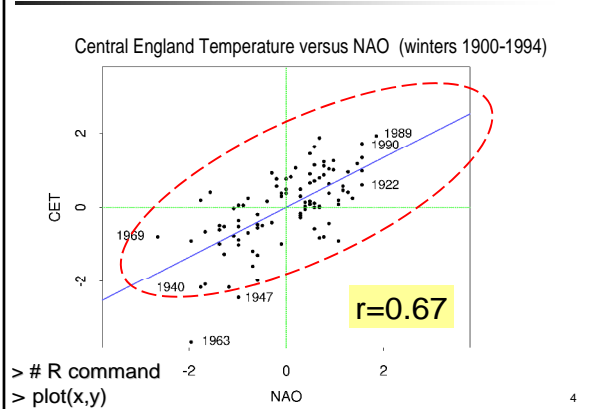
Note the "tied" values that occur in this small sample

```
> # R command  
> stripchart(x)
```

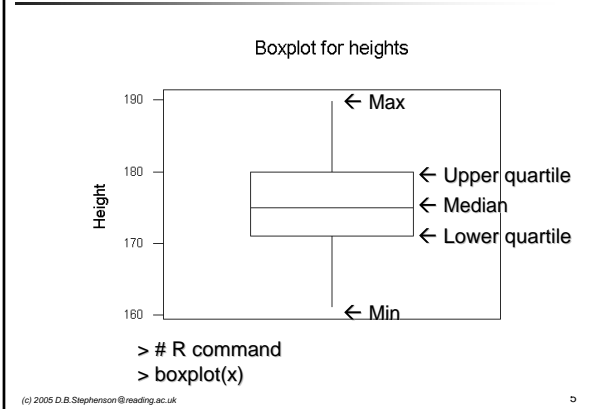
(c) 2005 D.B.Stephenson@reading.ac.uk

3

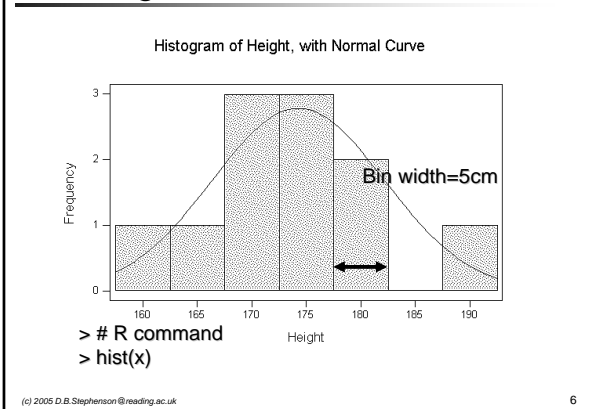
2. Scatter plot (2-d dotplot)



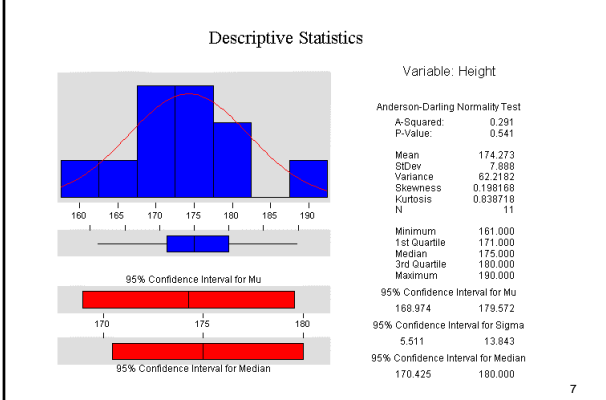
2. Boxplot



2. Histogram



2. Descriptive summary



2. Summary measures

- Centre/Location
 - Mean (m or x)
 - Median (x_{0.5})
- Scale/Spread
 - Standard deviation (s)
 - Interquartile range (IQR)
- Shape
 - Skewness (e.g. b1)
 - Kurtosis (e.g. b2)

2.3 The sample mean

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

Standard measure of the central location of the data.

> # R command
 > mean(x)

2.3 The standard deviation

$$s = \sqrt{(x - \bar{x})^2} = \sqrt{x^2 - (\bar{x})^2}$$

Std. Deviation=root mean squared deviation
Standard measure of the spread/scale of the data.

```
> # R command  
> sd(x)
```

(c) 2005 D.B.Stephenson@reading.ac.uk

10

2.3 Higher moments about mean

$$m_r = \overline{(x - \bar{x})^r}$$

Give information about the shape of the distribution
e.g. all odd moments are zero for a symmetric distribution

```
> # R command to do m4  
> m4<-mean((x-mean(x))^4)  
> m4
```

(c) 2005 D.B.Stephenson@reading.ac.uk

11

2.3 Skewness and kurtosis

$$\text{Skewness} = b_1 = m_3 / s^3$$

$$\text{Kurtosis} = b_2 = m_4 / s^4$$

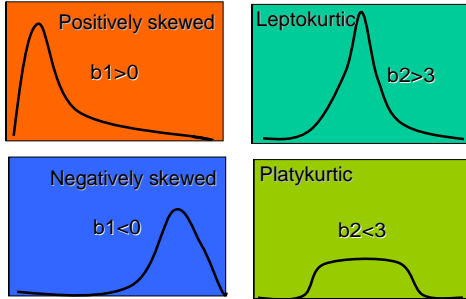
For normal (Gaussian) distribution:
Skewness=0 (symmetric) Kurtosis=3
Kurtosis > 3 "leptokurtic" (fat tails and sharp peak)
Kurtosis < 3 "platykurtic" (thin tails and flatter peak)

```
> # R commands  
> b1=mean(scale(x)^3)  
> b2=mean(scale(x)^4)
```

(c) 2005 D.B.Stephenson@reading.ac.uk

12

2.3 Shapes of distributions



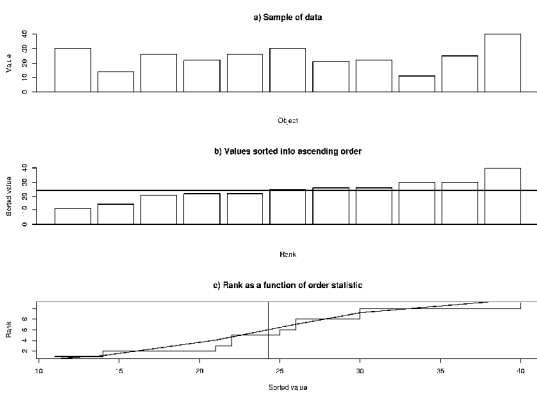
(c) 2005 D.B. Stephenson@reading.ac.uk

13

Resistant and Robust statistics

- **Resistant statistic** – one that is not overly sensitive to small or large outlier data e.g. IQR compared to max-min range.
- **Robust statistic** – one that is not dependent on the details of the probability distribution e.g. rank-statistics (median etc.)

2.4 Empirical cumulative distribution



Summary

- EDA
 - *know exactly how the data was produced*
 - *try to get hold of the raw untreated data*
 - *let the data speak for themselves*
- Summarise location, scale, shape
- Investigate outliers, tied values, and any other strange features
- Transform the data if necessary
