## Last time    VISUAL DATA DISPLAYS

In particular:    HISTOGRAMS based on
Frequency Distributions



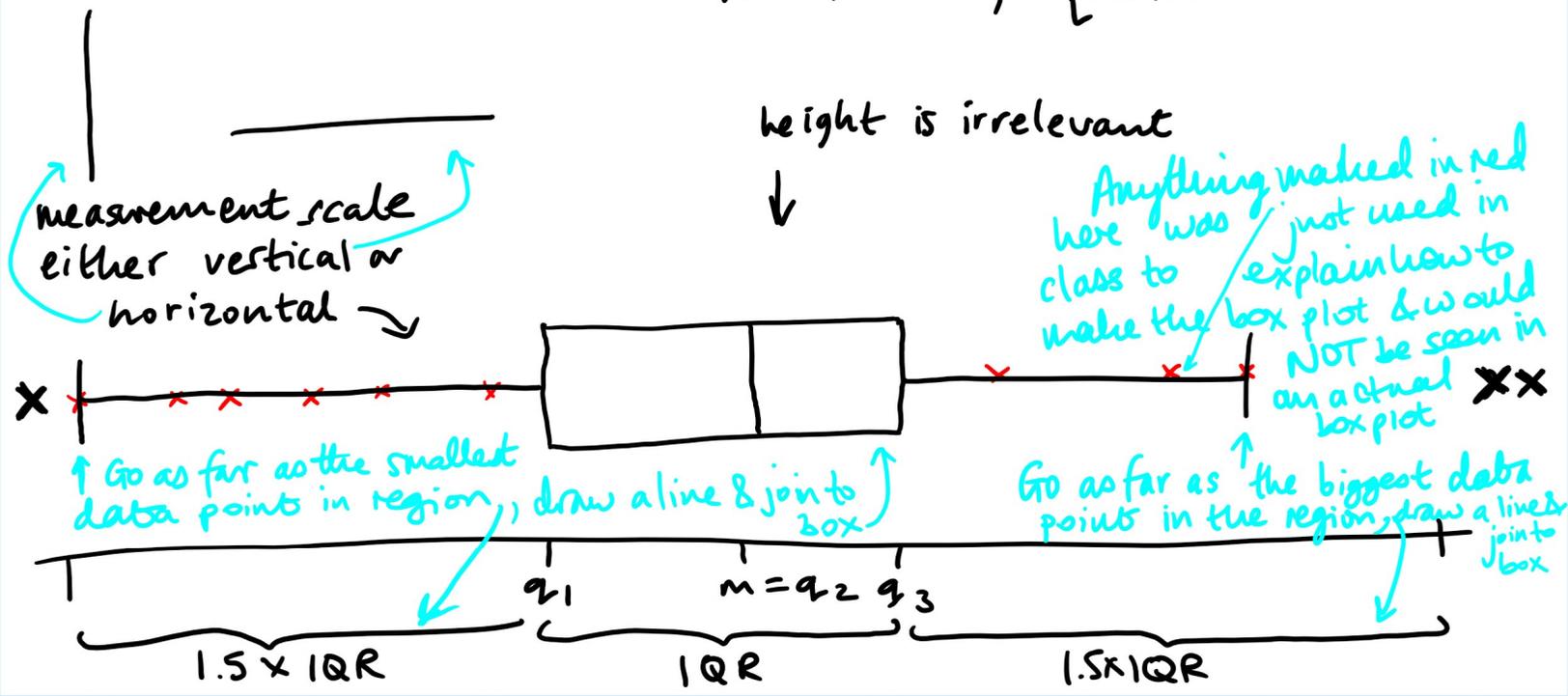↑ height of blocks = (relative) frequency in
corresponding bin

Range divided into
(equal-sized) "bins"

↑
If not equal,
then height = $\dfrac{\text{frequency}}{\text{width}}$

ALSO $q_1$= 1st quartile: 25% of data points below
$m = q_2 =$ 2nd quartile: 50% of data points below
$q_3 =$ 3rd quartile: 75% of data points below.

## 6.4  Box (and Whisker) Plots

— combines different features of data set (sample)
into one graph    ↓
min/max, quartiles

height is irrelevant
↓

measurement scale
either vertical or
horizontal ⤳

Anything marked in red
here was just used in
class to explain how to
make the box plot & would
NOT be seen in
an actual
box plot



↑ Go as far as the smallest
data points in region, draw a line & join to box

Go as far as the biggest data
points in the region, draw a line &
join to box

$q_1$        $m = q_2$ $q_3$

1.5 × IQR        IQR        1.5 × IQR

IQR = Interquartile Range = $q_3 - q_1$

In region either side of the box $\left(\text{from } q_1 \text{ to } q_3\right)$ draw "whiskers" out to most extreme data points lying in the region — these regions are $1.5 \times IQR$ either side

Then mark in all data points <u>outside</u> these two regions — these are called <u><u>outliers</u></u>

<span style="color:cyan">Definition of outlier</span> : ← <span style="color:cyan">data points outside of $1.5 \times IQR$ from box</span>

(1.5 — so that ~ 1% data points are outliers)

If a data point is $> 3 \times IQR$ outside the box, it is called an <u>extreme outlier</u>.

---

Can compare data sets (samples) by side-by-side box plots.

<u>Data Set 1</u>   $n = 12$.

10, 11, <u>16, 19</u>, 23, <u>31</u>, 33, 39, <u>50, 51</u>, 72, 105

$q_1 = 17.5$        $m = \dfrac{31+33}{2}$        $q_3 = 50.5$

<span style="color:cyan">$17.5 - 49.5 = -32$</span>        $= 32 = q_2$        $IQR = 50.5 - 17.5$
<span style="color:cyan">$50.5 + 49.5 = 100$</span>                   $= \boxed{33}$ ⭕

$1.5 \times IQR = 1.5 \times 33 = \boxed{49.5}$ ⭕
$3 \times IQR = 99$

---

<u>Data Set 2</u>   $n = 11$     1, 10, 23, 25, 27, 29, 30, 35, 36, 50, 89

<span style="color:green">$23 - 19.5 = 3.5$</span> $q_1$                     $m$"                        $q_3$"

$IQR = 36 - 23 = 13$, $1.5 \times IQR = 19.5$,        $3 \times IQR = 39$        <span style="color:green">$36 + 19.5 = 55.5$</span>

1.5×IQR: 10 down to -32

10  17.5  32  50.5  72  1.5×IQR: up to 100  105

10  17.5 ... 50.5 ... 72 ... X

1  10  23  29  36  89

3×IQR: up to 39+36=75

1.5×IQR: down to 3.5  1.5×IQR: up to 55.5

Outliers:     Data Set 1 :  105

Data Set 2 :  1, 89 ← extreme outlier

(36+33<89)

---

# 6.7 Probability Plots

— indicator of underlying probability distribution

↓

Histograms help indicate underlying prob. distr. but only reliable for large sample size $n$

---

Idea : We hypothesize the prob. distr. with pdf $f(x)$
(and this gives us cdf $F(x)$)

We check our guess with a probability plot:

How to make :     Sample is $\{x_1, \ldots, x_n\}$

→ Rank sample & rename so we have

$$x_{(1)} \leq x_{(2)} \leq \ldots \leq x_{(n)}$$

So $x_{(j)}$ is the $j$th sample point in numerical order

Idea is that $x_{(j)}$ should approximate the

$(100)\left(\dfrac{j}{n}\right)$ th percentile

So find, for each $j$, the number $y_j$ so that

$$\underset{\underset{\substack{\text{correction}\\ \text{factor as sample is}\\ \text{finite}}}{\uparrow}}{\dfrac{j-0.5}{n}} = \underset{\underset{\text{underlying population r.v.}}{\uparrow}}{P(X \le y_j)} = F(y_j) \;\leftarrow\; \substack{\text{the}\\ \text{hypothesized}\\ F}$$

Then plot $y_j$s against $x_{(j)}$s

— if guess is a good one then we get a straight line



increasing $y_j$s

Subjective judgment as to whether or not these points lie on a straight line.

increasing data points $(x_{(j)}s)$