## Review

We had two ways of describing the center and spread of a distribution

① $\bar{X}, S$

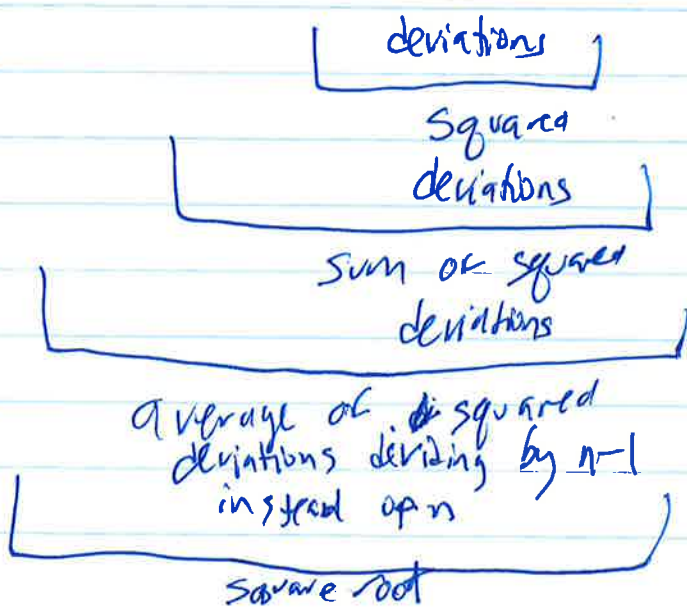mean, standard deviation

if either are present not as useful $\begin{cases} \text{sensitive to outliers (both measures)} \\ \text{doesn't show skewness in distribution} \\ \quad \text{eg. skewed to left} \\ \quad \text{or skewed to right} \end{cases}$

$$\bar{X} = \frac{X_1 + X_2 + X_3 + \cdots + X_n}{n} \qquad \frac{\text{sum of obs}}{\text{num of obs}}$$

$$= \frac{1}{n} \sum X_i \qquad \text{Sigma-notation}$$

$$S = \sqrt{\frac{1}{n-1} \sum (X_i - \bar{X})^2}$$

deviations

Squared deviations

Sum of squared deviations

average of squared deviations dividing by n-1 instead of n

square root

$S^2$ is variance

② min, $Q_1$, M, $Q_3$, max

a.k.a 5 number summary

minimum, 1st Quartile, Median, 3rd Quartile, maximum
0th percentile, 25th percentile, 50th percentile, 75th per, 100th

resistant to outliers (all measures, except min and max)

→ but in a modified box-plot you exclude outliers (or suspected outliers) from min and max in that case the min and max ~~would~~ as displayed would be resistant (to outliers)

Median — point where half the obs are above half are below

$Q_1$ — Median of ~~upper~~ lower half of obs
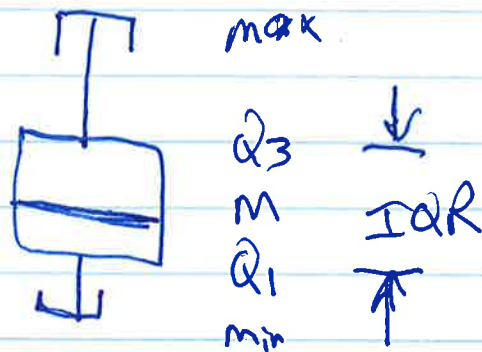
$Q_3$ — Median of upper half of obs

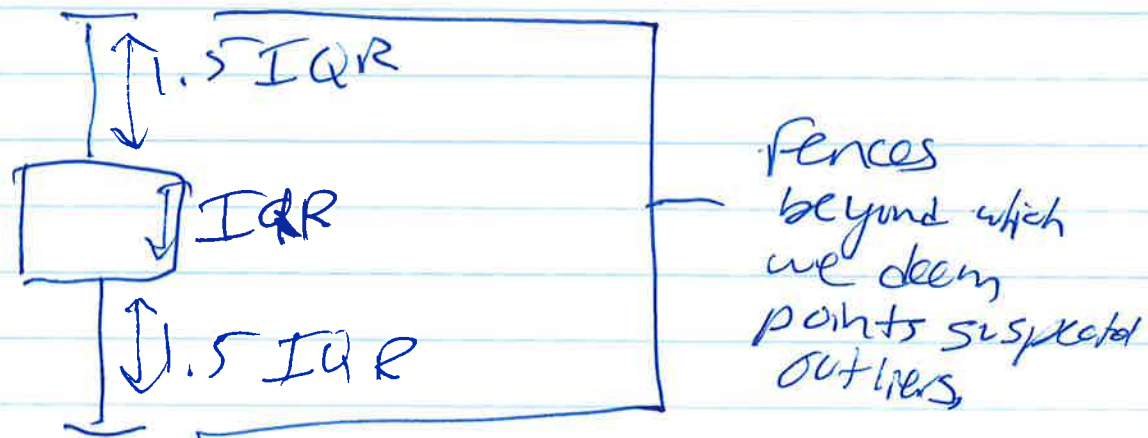min/max the smallest and largest obs (in a modified box plot you exclude suspect outliers from the min and max)

- The five number summary shows skewness of a distribution.

- The five number summary is displayed with a box plot

max

Q_3

M      IQR

Q_1

min

$$Q_3 - Q_1 = IQR = \text{Inner-quartile range.}$$

The IQR is useful for flagging suspected outliers

1.5 IQR

IQR

1.5 IQR

fences beyond which we deem points suspected outliers.

In a modified box plot the plotted min and max are inside fences and obs outside are plotted as asterisks.

Show homework ~~tb~~ 4

New: A _transformation_ is a function that transforms an old variable into a new ~~x~~ variable

$$X_{new} = F(X_{old})$$

Transformation is just another name for a function.

When we use the word transformation we think of the function as transforming the old variable

Examples

(a) IF $X_{old}$ is distance measured in kilometers and $X_{new}$ is ~~miles~~ distance measured in miles

$$X_{new} = .62 \, X_{old}$$

In familiar function notation
$$y = .62 X = F(x)$$

Example 2

$X_{old}$ = temperature measured in deg $f°$

$X_{new}$ = temperature measured in deg $c°$
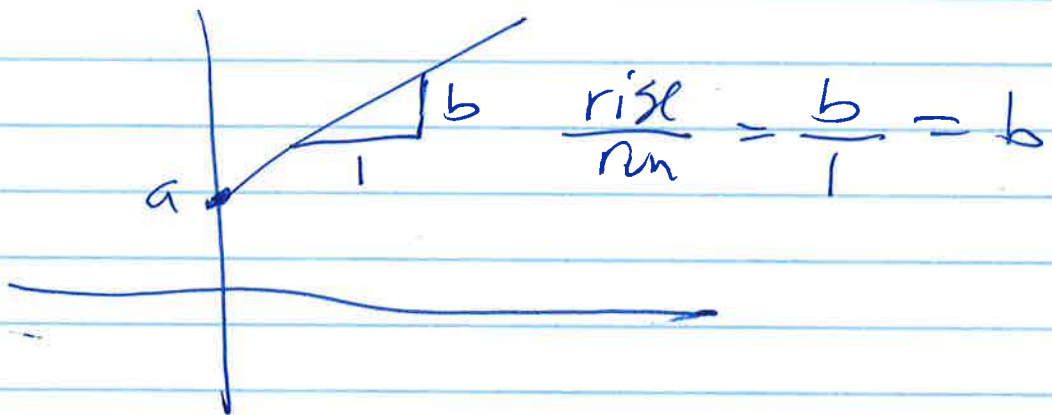
$$X_{new} = \frac{5}{9}(X_{old} - 32)$$
$$y = \frac{5}{9}x - 32$$

Unit conversions are usually
linear functions (graphs are lines)
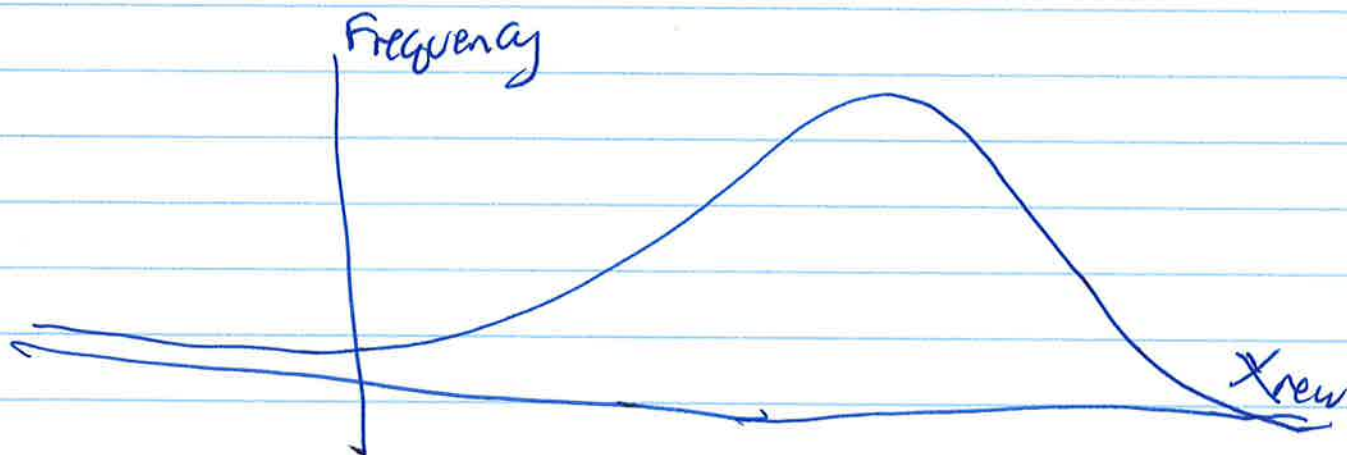
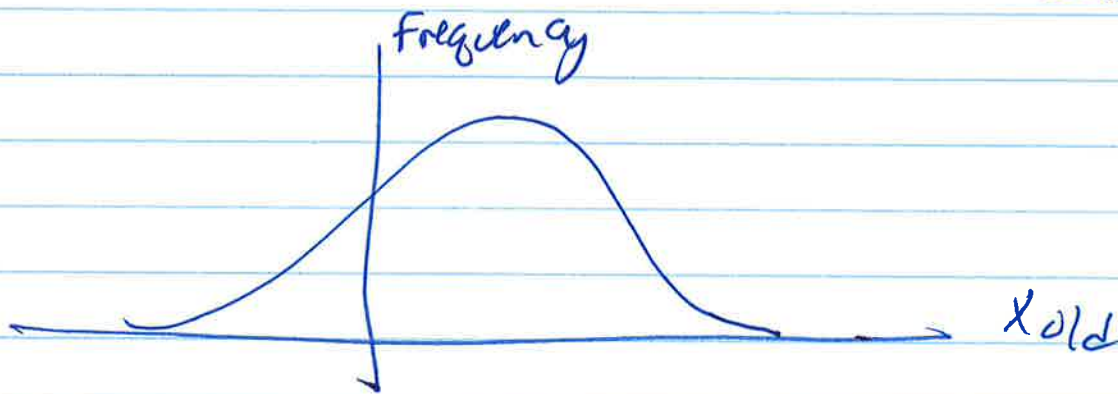These are called | linear transformations |

$$y = mx + b \quad \text{in familar notation}$$

$$X_{new} = a + b X_{old} \quad \text{book's notation}$$

$$\frac{rise}{run} = \frac{b}{1} = b$$

Linear transformations do not change the shape of a distribution

$X_{old}$ right skewed $\Leftrightarrow$ $X_{new}$ right skewed

$X_{old}$ symmetric unimodal $\Leftrightarrow$ $X_{new}$ symmetric unimodal

Frequency

$X_{old}$

Frequency

$X_{new}$

$$X_{new} = a + b X_{old}$$

a shifts histogram

b shrinks histogram if $|b| < 1$

expands " " $|b| > 1$

flips if $b < 0$

peaks, gaps, and skewness remain.

## Density Curves

A density curve is a smooth approximation to the irregular bars of a histogram.

For a histogram there were three plotting alternatives

A: Frequency — count of obs in bin
B: Relative Frequency — Freq / total number of obs
C: Density — Relative Freq / bin width

Why pick density? Suppose you collect more and more data. You are likely to see less and less noise in histogram. With A you will get more and more obs in each bin. Therefore, the vertical scale will keep increasing. Not so with B and C.

Another thing to do is to decrease the width of bin as you increase the number of data points. This will give a smoother and smoother histogram.

B will change (vertical scale) as you decrease bin width C wont.

Thats why we like C.

As we approach a histogram
with more and more data and
smaller and smaller bin width
we approach a histogram

[Show with stat crunch]

For A (Frequency) the sum of bin
heights is n

For B ~~xxxxx~~ the sum of bin heights is 1.

For C, the area of the bins is 1

For density curves, like for C.

(1) the area under curve is 1
(2) the curve is always on or above
the horizontal axis.

A density curve is any curve that satisfies
(1) and (2). A distribution is completely
described by its density curve.

So instead of giving the 5-number summary of a variable it would be better to give the density curve.

Problem though: You can never have enough data to be sure you know the density curve.

There are ways to estimate density curves from data. But StatCrunch does not do this

What StatCrunch will do instead is overlay your histogram of data with a density curve of a standard distribution

There are infinitely many distributions (any density curve determines one)

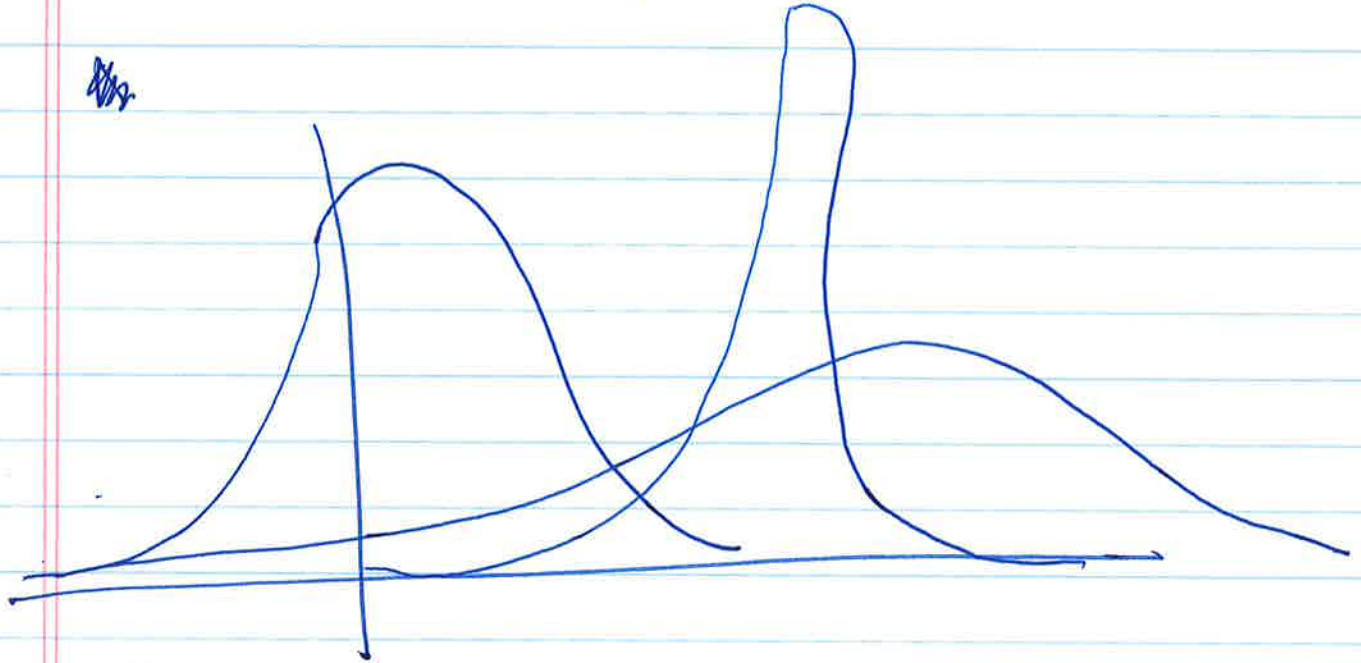Only a few have names (the standard distributions)

The most important standard distribution is the <u>Normal Distribution</u>

Its density curve is the <u>bell curve</u>.

Bell Curves are determined by
their mean and standard deviation

[ Mean of data written $\bar{X}$
  Standard deviation of data written $S$

[ Mean of density curve written $\mu$
  Std dev $\sigma$

We'll give precise defns next week