# Measures of Central Tendency:

## Mean, Median, Mode

# Measures of Central Tendency

- **Measure of central tendency provides a very convenient way of describing a set of scores with a single number that describes the PERFORMANCE of the group.**

- **It is also defined as a single value that is used to describe the "center" of the data.**

- **There are three commonly used measures of central tendency. These are the following:**

    - **MEAN**

    - **MEDIAN**

    - **MODE**

# MEAN

- **It is the most commonly used measure of the center of data**

- **It is also referred as the "arithmetic average"**

  - **Computation of  Sample Mean**

$$\overline{X} = \frac{\Sigma X}{N} = \frac{x_1 + x_2 + x_3 + \dots x_n}{N}$$

  - **Computation of  the Mean for Ungrouped Data**

$$\overline{X} = \frac{\Sigma x}{n} \qquad \overline{X} = \frac{\Sigma f x}{n}$$

# MEAN

**Example:**

**Scores of 15 students in Mathematics I quiz consist of 25 items. The highest score is 25 and the lowest score is 10. Here are the scores: 25, 20, 18, 18, 17, 15, 15, 15, 14, 14, 13, 12, 12, 10, 10. Find the mean in the following scores.**

**x (scores)**

| | |
|---|---|
| 25 | 14 |
| 20 | 14 |
| 18 | 13 |
| 18 | 12 |
| 17 | 12 |
| 15 | 10 |
| 15 | 10 |
| 15 | |

$$\overline{X} = \frac{\Sigma x}{n}$$

$$= \frac{228}{15}$$

$$= 15.2$$

# MEAN

$$\overline{X} = 15.2$$

**Analysis:**

The average performance of 15 students who participated in mathematics quiz consisting of 25 items is 15.20. The implication of this is that student who got scores below 15.2 did not perform well in the said examination. Students who got scores higher than 15.2 performed well in the examination compared to the performance of the whole class.

# MEAN

**Example:**

Find the Grade Point Average (GPA) of Paolo Adade for the first semester of the school year 2013-2014. Use the table below:

| Subjects | Grade ($X_i$) | Units ($w_i$) | ($X_i$) ($w_i$) |
|---|---|---|---|
| **BM 112** | 1.25 | 3 | 3.75 |
| BM 101 | 1.00 | 3 | 3.00 |
| AC 103 | 1.25 | 6 | 7.50 |
| EC 111 | 1.00 | 3 | 3.00 |
| MG 101 | 1.50 | 3 | 4.50 |
| MK 101 | 1.25 | 3 | 3.75 |
| FM 111 | 1.50 | 3 | 4.50 |
| PE 2 | 1.00 | 2 | 2.00 |
| | | $\Sigma(w_i) = 26$ | $\Sigma(X_i)w_i) = 32.00$ |

# MEAN

$$\overline{X} = \frac{\Sigma(X_i)w_i}{\Sigma(w_i)}$$

$$= \frac{32}{26}$$

$$= 1.23$$

The Grade Point Average of Paolo Adade for the
first semester SY 2013-2014
Is 1.23.

# MEAN

**Mean for Grouped Data**

*Grouped data* are the data or scores that are arranged in a frequency distribution.

*Frequency distribution* is the arrangement of scores according to category of classes including the frequency.

*Frequency* is the number of observations falling in a category.

# MEAN

**The only one formula in solving the mean for grouped data is called *midpoint method.* The formula is:**

$$\overline{X} = \frac{\Sigma f\, x_m}{n}$$

Where $\overline{X}$ = mean value

$x_m$ = midpoint of each class or category

f = frequency in each class or category

$\Sigma f\, x_m$ = summation of the product of $f\, x_m$

# MEAN

**Steps in Solving Mean for Grouped Data**

1. **Find the midpoint or class mark ($Xm$) of each class or category using the formula**  $Xm = \dfrac{LL + LU}{2}$ .

2. Multiply the frequency and the corresponding class mark $f\, x_m$.

3. **Find the sum of the results in step 2.**

4. Solve the mean using the formula

$$\bar{X} = \frac{\Sigma\, f\, x_m}{n}$$

# MEAN

**Example:**

Scores of 40 students in a science class consist of 60 items and they are tabulated below.

| X | f | Xm | fXm |
|---|---|----|-----|
| 10 – 14 | 5 | 12 | 60 |
| 15 – 19 | 2 | 17 | 34 |
| 20 – 24 | 3 | 22 | 66 |
| 25 – 29 | 5 | 27 | 135 |
| 30 – 34 | 2 | 32 | 64 |
| 35 – 39 | 9 | 37 | 333 |
| 40 – 44 | 6 | 42 | 252 |
| 45 – 49 | 3 | 47 | 141 |
| 50 - 54 | 5 | 52 | 260 |
|  | n = 40 |  | Σ f Xm = 1 345 |

$$\bar{X} = \frac{\Sigma f X_m}{n}$$

$$= \frac{1\,345}{40}$$

$$= 33.63$$

# MEAN

**Analysis:**

**The mean performance of 40 students in science quiz is 33.63. Those students who got scores below 33.63 did not perform well in the said examination while those students who got scores above 33.63 performed well.**

# MEAN

**Properties of the Mean**

- **It measures <span style="color:orange">stability</span>. Mean is the most stable among other measures of central tendency because every score contributes to the value of the mean.**

- **The sum of each score's distance from the mean is zero.**

- **It may easily affected by the extreme scores.**

- **It can be applied to interval level of measurement.**

- **It may not be an actual score in the distribution.**

- **It is very easy to compute.**

# MEAN

**When to Use the Mean**

- Sampling stability is desired.

- Other measures are to be computed such as standard deviation, coefficient of variation and skewness.

# MEDIAN

- Median is what divides the scores in the distribution into two equal parts.

- Fifty percent (50%) lies below the median value and 50% lies above the median value.

- It is also known as the middle score or the 50th percentile.

# MEDIAN

**Median of Ungrouped Data**

1.  Arrange the scores (from lowest to highest or highest to lowest).

2.  Determine the middle most score in a distribution if *n* is an *odd number* and get the *average* of the two middle most scores if *n* is an *even number.*

Example 1: Find the median score of 7 students in an English class.

| x (score) |
|:---:|
| 19 |
| 17 |
| 16 |
| **15** |
| 10 |
| 5 |
| 2 |

# MEDIAN

**Example: Find the median score of 8 students in an English class.**

$$x\ (score)$$

30

19

17

*16*

*15*

10

5

2

$$\tilde{x} = \frac{16 + 15}{2}$$

$$\tilde{x} = 15.5$$

# MEDIAN

**Median of Grouped Data**

**Formula:**

$$\tilde{x} = L_B + \frac{\frac{n}{2} - cfp}{fm} \times c.i$$

$\tilde{X}$ = median value

MC = median class is a category containing the $\frac{n}{2}$

$L_B$ = lower boundary of the median class (MC)

cfp = cumulative frequency before the median class if the scores are arranged from lowest to highest value

fm = frequency of the median class

c.i = size of the class interval

# MEDIAN

**Steps in Solving Median for Grouped Data**

1. Complete the table for cf<.

2. Get $\dfrac{n}{2}$ of the scores in the distribution so that you can identify MC.

3. Determine $L_B$, cfp, fm, and c.i.

4. Solve the median using the formula.

# MEDIAN

Example: Scores of 40 students in a science class consist of 60 items and they are tabulated below. The highest score is 54 and the lowest score is 10.

| X | f | cf< |
|---|---|---|
| 10 – 14 | 5 | 5 |
| 15 – 19 | 2 | 7 |
| 20 – 24 | 3 | 10 |
| 25 – 29 | 5 | 15 |
| 30 – 34 | 2 | **17 (cfp)** |
| 35 – 39 | **9 (fm)** | 26 |
| 40 – 44 | 6 | 32 |
| 45 – 49 | 3 | 35 |
| 50 – 54 | 5 | 40 |
| | n = 40 | |

# MEDIAN

**Solution:**

$$\frac{n}{2} = \frac{40}{2} = 20$$

The category containing $\frac{n}{2}$ is 35 –39.

LL of the MC = 35

$L_n$ = 34.5

cfp = 17

fm = 9

c.i = 5

$$\tilde{x} = L_B + \frac{\frac{n}{2} - cfp}{fm} \times c.i$$

$$= 34.5 + \frac{20 - 17}{9} \times 5$$

$$= 34.5 + 15/9$$

$$\tilde{x} = 36.17$$

# MEDIAN

**Properties of the Median**

- It may not be an actual observation in the data set.

- It can be applied in ordinal level.

- It is not affected by extreme values because median is a positional measure.

**When to Use the Median**

- The exact midpoint of the score distribution is desired.

- There are extreme scores in the distribution.

# MODE

The *mode* or the *modal score* is a score or scores that occurred most in the distribution.

It is classified as unimodal, bimodal, trimodal or mulitimodal.

*Unimodal* is a distribution of scores that consists of only one mode.

*Bimodal* is a distribution of scores that consists of two modes.

*Trimodal* is a distribution of scores that consists of three modes or *multimodal* is a distribution of scores that consists of more than two modes.

# MODE

Example: Scores of 10 students in Section A, Section B and Section C.

| Scores of Section A | Scores of Section B | Scores of Section C |
|---|---|---|
| 25 | 25 | 25 |
| 24 | 24 | 25 |
| 24 | 24 | 25 |
| 20 | 20 | 22 |
| 20 | 18 | 21 |
| 20 | 18 | 21 |
| 16 | 17 | 21 |
| 12 | 10 | 18 |
| 10 | 9 | 18 |
| 7 | 7 | 18 |

# MODE

The score that appeared most in Section A is 20, hence, the mode of Section A is *20.* There is only one mode, therefore, score distribution is called *unimodal.*

The modes of Section B are *18* and *24,* since both 18 and 24 appeared twice. There are two modes in Section B, hence, the distribution is a *bimodal distribution.*

The modes for Section C are *18, 21,* and *25.* There are three modes for Section C, therefore, it is called a *trimodal* or *multimodal distribution.*

# MODE

**Mode for Grouped Data**

**In solving the mode value in grouped data, use the formula:**

$$\hat{X} = L_B + \frac{d_1}{d_1 + d_2} \times c.i$$

$L_B$ = lower boundary of the modal class

**Modal Class (MC)** = is a category containing the highest frequency

$d_1$ = difference between the frequency of the modal class and the frequency above it, when the scores are arranged from lowest to highest.

$d_2$ = difference between the frequency of the modal class and the frequency below it, when the scores are arranged from lowest to highest.

$c.i$ = size of the class interval

# MODE

Example: Scores of 40 students in a science class consist of 60 items and they are tabulated below.

| x | f |
|---|---|
| 10 - 14 | 5 |
| 15 - 19 | 2 |
| 20 - 24 | 3 |
| 25 - 29 | 5 |
| 30 - 34 | 2 |
| 35 - 39 | 9 |
| 40 - 44 | 6 |
| 45 - 49 | 3 |
| 50 - 54 | 5 |
|  | n = 40 |

# MODE

Modal Class = 35 – 39

LL of MC = 35

$L_B$ = 34.5

$d_1$ = 9 – 2 = 7

$d_2$ = 9 – 6 = 3

c.i = 5

$$\hat{X} = L_B + \frac{d_1}{d_1 + d_2} \times c.i$$

$$= 34.5 + \frac{7}{7 + 3} \times 5$$

= 34. 5 + 35/10

$\hat{X}$ = 38

The mode of the score distribution that consists of 40 students is 38, because 38 occurred several times.

# MODE

## Properties of the Mode

- It can be used when the data are qualitative as well as quantitative.

- It may not be unique.

- It is affected by extreme values.

- It may not exist.

## When to Use the Mode

- When the "typical" value is desired.
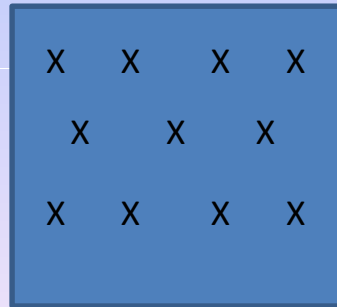
- When the data set is measured on a nominal scale.

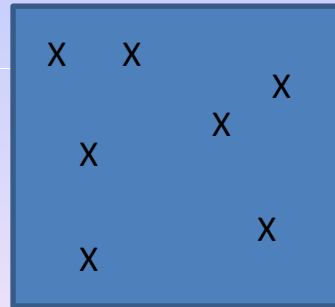# Nearest Neighbour Analysis

# What is it?

- Indicates spatial distribution of area-from average distance between each point and nearest neighbour

- 3 types of pattern:

- **Regular**

- **Clustered**

- **Random**



NNI=0.0          NNI=2.15          NNI=1.0

# Formula

- NNI=2$\bar{D}$ $\sqrt{(N/A)}$

- $\bar{D}$: average distance between each point & its nearest neighbour (**∑d/N**) [**d**=each individual distance]

- **N**: number of studied points

- **A**: size of the studied area ! SAME UNITS!

# Example 1-flat land

| Settlement | Nearest Neighbour | Distance (km) |
|---|---|---|
| Cierny Majer | Kosuty | 1.6 |
| Kosuty | Cierny majer | 1.6 |
| Hed | Cierny Brod | 0.9 |
| Cierny Brod | Hed | 0.9 |
| Mostova | Cierny Brod | 2 |
| Cierna Voda | Cierny Brod | 2.4 |
| Stary Haj | Vozokany | 3 |
| Vozokany | Stary Haj | 3 |
| Degessky majer | Dolna luka | 3.4 |
| Dolna luka | Matuskovo | 2.5 |
| Matuskovo | Budic | 2.2 |
| Budic | Matuskovo | 2.2 |
| | | =25.7 |

# Example 1-flat land

- Total=25.7
- Area= 100km2
- Villages_=12
- NNI=2D̄ $\sqrt{(N/A)}$
- 2D̄ :25.7/12=2.14*2=4.28
- NNI=4.28$\sqrt{(12/100)}$
- =1.48
- $\rightarrow$**regular**
- A lot of free land on flat plains

# Example 2-hilly land

| Settlement | Nearest Neighbour | Distance (km) |
|---|---|---|
| Mlyniste | Slace | 1.2 |
| Slace | Mlyniste | 1.2 |
| Prostredny vrch | Skycov | 1.1 |
| Skycov | Prostredny vrch | 1.1 |
| Piesky | Dlhe diely | 0.9 |
| Dlhe diely | Piesky | 0.9 |
| Breziny | Dlhe Diely | 1.8 |
| Kosiar | Hostianske | 0.6 |
| Hostianske | Kosiar | 0.6 |
| Hrusov | Rybnik | 1 |
| Rybnik | Rybniky | 0.9 |
| Borinky | Modos | 1.1 |
| Modos | Masirov stal | 0.8 |
| Uholna bana | Modos | 0.9 |
| Masirov stal | Modos | 0.8 |
| Jedlove Kostolany | Lukacov stal | 1 |
| Lukacov stal | Levasovsky stal | 0.4 |
| Levasovksy stal | Lukacov stal | 0.4 |
| Borisko | Levasovsky stal | 0.7 |
| Koborno | Kraje | 3.2 |
| Kraje | Klizske Lucky | 1.5 |
| Klizske Lucky | Kraje | 1.5 |
| Drahozicka huta | Klizske Lucky | 2.5 |

# Example 2-hilly land

- NNI=1.12 $\sqrt{(23/100)}$

-  NNI=0.536

- $\rightarrow$clustered $\rightarrow$regular


- Valleys- ecumene settlements

# Drawbacks

- 'straight line assumption': on a map the distances through mountains not considered

- Where do you measure from? How do you determine the centre of the settlement exactly?

- Measure distance by road/straight line?

- What settlements to include? Size limit?

- Effect of paired distributions?

- Controlling factors e.g. Soil type, relief

# Practical use?

Comparison of distributions-
quantifiable measure of a
distribution pattern

# Standard Statistical Distributions Normal, Poisson, Binomial

Normal distribution describes continuous data which have a symmetric distribution, with a characteristic 'bell' shape.
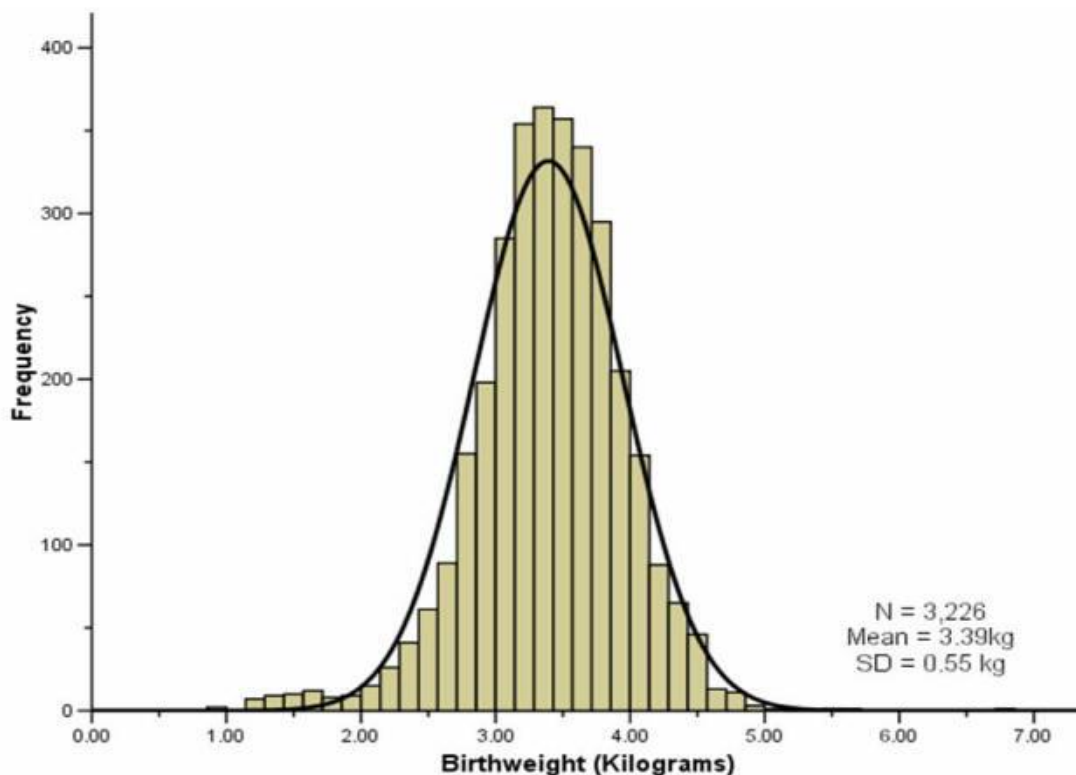
Binomial distribution describes the distribution of binary data from a finite sample. Thus it gives the probability of getting r events out of n trials.

Poisson distribution describes the distribution of binary data from an infinite sample. Thus it gives the probability of getting r events in a population.

**The Normal Distribution**

It is often the case with medical data that the histogram of a continuous variable obtained from a single measurement on different subjects will have a characteristic `bell-shaped' distribution known as a Normal distribution. One such example is the histogram of the birth weight (in kilograms) of the 3,226 new born babies shown in Figure 1.

*Figure 1 Distribution of birth weight in 3,226 newborn babies (data from O' Cathain et al 2002)*

To distinguish the use of the same word in normal range and Normal distribution we have used a lower and upper case convention throughout.

The histogram of the sample data is an estimate of the population distribution of birth weights in new born babies. This population distribution can be estimated by the superimposed smooth `bell-shaped' curve or `Normal' distribution shown. We presume that if we were able to look at the entire population of new born babies then the distribution of birth weight would have exactly the Normal shape. We often infer, from a sample whose histogram has the approximate Normal shape, that the population will have exactly, or as near as makes no practical difference, that Normal shape.

The Normal distribution is completely described by two parameters μ and σ, where μ represents the population mean, or centre of the distribution, and σ the population standard deviation. It is symmetrically distributed around the mean. Populations with small values of the standard deviation σ have a distribution concentrated close to the centre μ; those with large standard deviation have a distribution widely spread along the measurement axis. One mathematical property of the Normal distribution is that exactly 95% of the distribution lies between

$$\mu-(1.96 \times \sigma) \text{ and } \mu+(1.96 \times \sigma)$$

Changing the multiplier 1.96 to 2.58, exactly 99% of the Normal distribution lies in the corresponding interval.

In practice the two parameters of the Normal distribution, μ and σ, must be estimated from the sample data. For this purpose a random sample from the population is first taken. The sample mean $\tilde{x}$ and the sample standard deviation, $SD(\bar{x}) = \sqrt{SSD(\bar{x})} = S$ , are then calculated. If a sample is taken from such a Normal distribution, and provided the sample is not too small, then approximately 95% of the sample lie within the interval:
$$\bar{x}-[1.96 \times SD(\bar{x})] \text{ to } \bar{x}+[1.96 \times SD(\bar{x})]$$

This is calculated by merely replacing the population parameters μ and σ by the sample estimates $\tilde{x}$ and s in the previous expression.

In appropriate circumstances this interval may estimate the reference interval for a particular laboratory test which is then used for diagnostic purposes.

We can use the fact that our sample birth weight data appear Normally distributed to calculate a reference range. We have already mentioned that about 95% of the observations (from a Normal distribution) lie within ±1.96 SDs of the mean. So a reference range for our sample of babies, using the values given in the histogram above, is:

3.39 - [1.96 x 0.55]  to  3.39 + [1.96 x 0.55]

2.31kg to 4.47kg

A baby's weight at birth is strongly associated with mortality risk during the first year and, to a lesser degree, with developmental problems in childhood and the risk of various diseases in adulthood. If the data are not Normally distributed then we can base the normal reference range on the observed percentiles of the sample, i.e. 95% of the observed data lie between the 2.5 and 97.5 percentiles. In this example, the percentile-based reference range for our sample was calculated as 2.19kg to 4.43kg.

Most reference ranges are based on samples larger than 3500 people. Over many years, and millions of births, the WHO has come up with a normal birth weight range for new born babies. These ranges represent results than are acceptable in newborn babies and actually cover the middle 80% of the population distribution, i.e. the 10th to 90th centiles. Low birth weight babies are usually defined (by the WHO) as weighing less than

2500g (the 10th centile) regardless of gestational age, and large birth weight babies are defined as weighing above 4000kg (the 90th centile). Hence the normal birth weight range is around 2.5kg to 4kg. For our sample data, the 10th to 90th centile range was similar, 2.75 to 4.03kg.


## The Binomial Distribution

If a group of patients is given a new drug for the relief of a particular condition, then the proportion $p$ being successively treated can be regarded as estimating the population treatment success rate $\pi$.

The sample proportion $p$ is analogous to the sample mean $\tilde{\bar{x}}$, in that if we score zero for those $s$ patients who fail on treatment, and 1 for those $r$ who succeed, then $p=r/n$, where $n=r+s$ is the total number of patients treated. Thus $p$ also represents a mean. Data which can take only a binary (0 or 1) response, such as treatment failure or treatment success, follow the binomial distribution provided the underlying population response rate does not change. The binomial probabilities are calculated from:

$$P(r \text{ responses out of } n) = \frac{n!}{r!(n-r)!} \pi^r (1-\pi)^{n-r}$$

…for successive values of R from 0 through to n. In the above, $n!$ is read as "n factorial" and $r!$ as "r factorial". For $r=4$, $r!=4\times3\times2\times1=24$. Both 0! and 1! are taken as equal to 1. The shaded area marked in Figure 2 (below) corresponds to the above expression for the binomial distribution calculated for each of $r=8,9,...,20$ and then added. This area totals 0.1018. So the probability of eight or more responses out of 20 is 0.1018.
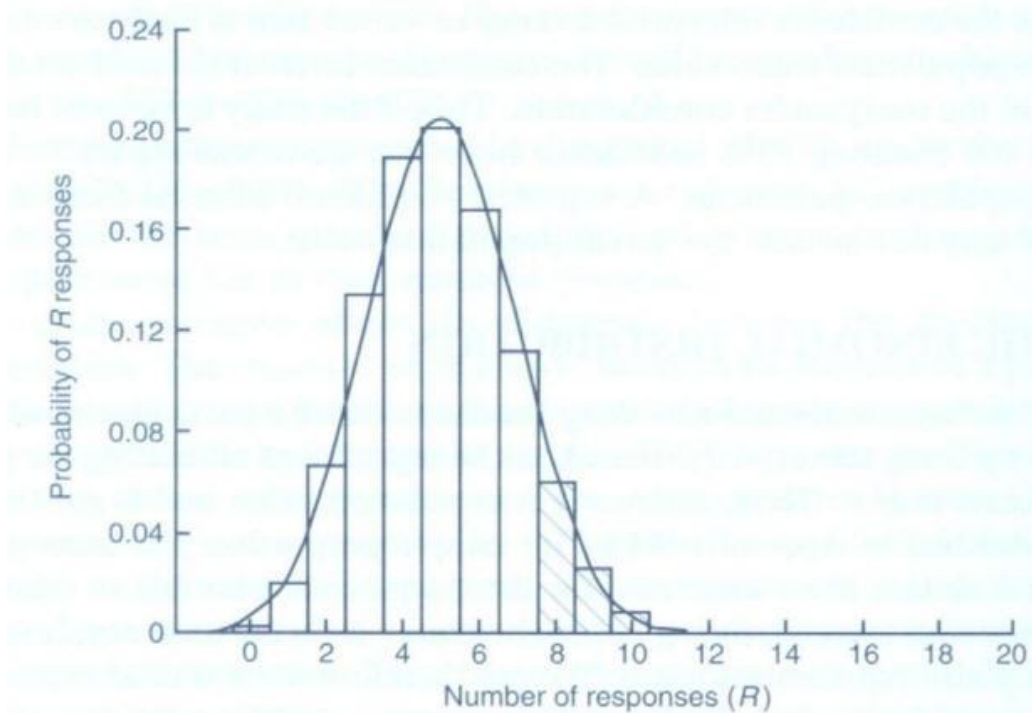
For a fixed sample size $n$ the shape of the binomial distribution depends only on $\pi$. Suppose $n = 20$ patients are to be treated, and it is known that on average a quarter, or $\pi =0.25$, will respond to this particular treatment. The number of responses actually observed can only take integer values between 0 (no responses) and 20 (all respond). The binomial distribution for this case is illustrated in Figure 2.

The distribution is not symmetric, it has a maximum at five responses and the height of the blocks corresponds to the probability of obtaining the particular number of responses from the 20 patients yet to be treated. It should be noted that the expected value for $r$, the number of successes yet to be observed if we treated $n$ patients, is ($n \times \pi$). The potential variation about this expectation is expressed by the corresponding standard deviation:

$$SD(r) = \sqrt{n\pi(1-\pi)}$$

Figure 2 also shows the Normal distribution arranged to have $\mu = n\pi = 5$ and $\sigma = \sqrt{[n\pi(1 - \pi)]} = 1.94$, superimposed on to a binomial distribution with $\pi = 0.25$ and n = 20. The Normal distribution describes fairly precisely the binomial distribution in this case.      If n is small, however, or $\pi$ close to 0 or 1, the disparity between the Normal and binomial distributions with the same mean and standard deviation increases and the Normal distribution can no longer be used to approximate the binomial distribution. In such cases the probabilities generated by the binomial distribution itself must be used. It is also only in situations in which reasonable agreement exists between the distributions that we would use the confidence interval expression given previously. For technical reasons, the expression given for a confidence interval for a proportion is an approximation. The approximation will usually be quite good provided $p$ is not too close to 0 or 1, situations in which either almost none or nearly all of the patients respond to treatment. The approximation improves with increasing sample size $n$.

*Figure 2: Binomial distribution for n=20 with $\pi$=0.25 and the Normal approximation*

## The Poisson Distribution

The Poisson distribution is used to describe discrete quantitative data such as counts in which the population size $n$ is large, the probability of an individual event $\pi$ is small, but the expected number of events, $n\pi$, is moderate (say five or more). Typical examples are the number of deaths in a town from a particular disease per day, or the number of admissions to a particular hospital.

*Example*
Wight et al (2004) looked at the variation in cadaveric heart beating organ donor rates in the UK. They found that there were 1330 organ donors, aged 15-69, across the UK for the two years 1999 and 2000 combined. Heart-beating donors are patients who are seriously ill in an intensive care unit (ICU) and are placed on a ventilator.

Now it is clear that the distribution of the number of donors takes integer values only, thus the distribution is similar in this respect to the binomial. However, there is no theoretical limit to the number of organ donors that could happen on a particular day. Here the population is the UK population aged 15-69, over two years, which is over 82 million person years, so in this case each member can be thought to have a very small probability of actually suffering an event, in this case being admitted to a hospital ICU and placed on a ventilator with a life threatening condition.

The mean number of organ donors per day over the two year period is calculated as:

$$r = \frac{1330}{(365+365)} = \frac{1330}{730} = 1.82 \text{ organ donations per day}$$

It should be noted that the expression for the mean is similar to that for $\pi$, except here multiple data values are common; and so instead of writing each as a distinct figure in the numerator they are first grouped and counted. For data arising from a Poisson distribution the standard error, that is the standard deviation of $r$, is estimated by $SE(r) = \sqrt{(r/n)}$, where $n$ is the total number of days (or an alternative time unit). Provided the organ donation rate is not too low, a 95% confidence interval for the underlying (true) organ donation rate $\lambda$ can be calculated in the usual way:

$$r-[1.96\times SE(r)]\text{to}r+[1.96\times SE(r)]r-[1.96\times SE(r)]\text{to}r+[1.96\times SE(r)]$$

In the above example $r=1.82$, $SE(r)=\sqrt{(1.82/730)}=0.05$, and therefore the 95% confidence interval for $\lambda$ is 1.72 to 1.92 organ donations per day. Exact confidence intervals can be calculated as described by Altman et al. (2000).

The Poisson probabilities are calculated from:

$$P(r \text{ responses})=\frac{\lambda^r}{r!}e^{-\lambda}P(r \text{ responses})=\frac{\lambda^r}{r!}e^{-\lambda}$$

…for successive values of $r$ from 0 to infinity. Here $e$ is the exponential constant 2.7182…, and $\lambda$ is the population rate which is estimated by $r$ in the example above.

*Example*

Suppose that before the study of Wight et al. (2004) was conducted it was expected that the number of organ donations per day was approximately two. Then assuming $\lambda = 2$, we would anticipate the probability of 0 organ donations in a given day to be $(2^0/0!)e^{-2} =e^{-2} = 0.135$. (Remember that $2^0$ and $0!$ are both equal to 1.) The probability of one organ donation would be $(2^1/1!)e^{-2} = 2(e^{-2}) = 0.271$. Similarly the probability of two organ donations per day is $(2^2/2!)e^{-2}= 2(e^{-2}) = 0.271$; and so on to give for three donations 0.180, four donations 0.090, five donations 0.036, six donations 0.012, etc. If the study is then to be conducted over 2 years (730 days), each of these probabilities is multiplied by 730 to give the expected number of days during which 0, 1, 2, 3, etc. donations will occur. These expectations are 98.8, 197.6, 197.6, 131.7, 26.3, 8.8 days. A comparison can then be made between what is expected and what is actually observed.
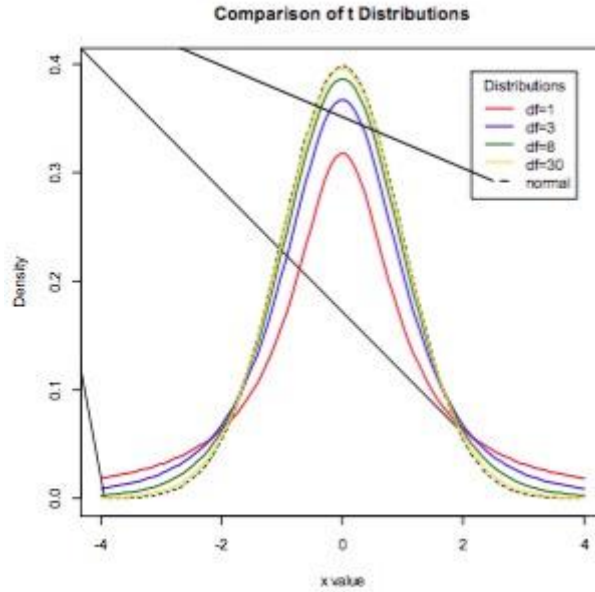
## Other Distributions

A brief description of some other distributions are given for completeness.

## t-distribution

Student's *t*-distribution is a continuous probability distribution with a similar shape to the Normal distribution but with wider tails. *t*-distributions are used to describe samples which have been drawn from a population, and the exact shape of the distribution varies with the sample size. The smaller the sample size, the more spread out the tails, and the larger the sample size, the closer the *t*-distribution is to the Normal distribution (Figure 3). Whilst in general the Normal distribution is used as an approximation when estimating means of samples from a Normally-distribution population, when the same size is small (say n<30), the *t*-distribution should be used in preference.

*Figure 3. The t-distribution for various sample sizes. As the sample size increases, the t-distribution more closely approximates the Normal.*

Comparison of t Distributions

## Chi-squared distribution

The chi-squared distribution is continuous probability distribution whose shape is defined by the number of degrees of freedom. It is a right-skew distribution, but as the number of degrees of freedom increases it approximates the Normal distribution (Figure 4). The chi-squared distribution is important for its use in chi-squared tests. These are often used to test deviations between observed and expected frequencies, or to determine the independence between categorical variables. When conducting a chi-squared test, the probability values derived from chi-squared distributions can be looked up in a statistical table.

*Figure 4. The chi-squared distribution for various degrees of freedom. The distribution becomes less right-skew as the number of degrees of freedom increases.*