

**Mr G's Little Book on**

**Probability**

**Distributions**

## Preface

MS Word supplies over 10 000 characters ranging from the sublime to the ridiculous but not those key expressions required in statistics.  $\bar{x}$  is in Arial Unicode and looks reasonable but capital  $x$  ie  $\bar{X}$  is undersized. In both cases Word then shifts open the line spacing.

Placing a circumflex above a Greek letter is all but impossible – the best I could realise was  $\sigma^{\wedge}$ .

Determined not to use equation editor a definite integral between  $a$  and  $b$  is  $\int_a^b$  and the result in square brackets is  $[ \text{some } f(x) ]_a^b$ . Where the limits have suffixes the result looks a bit clumsy  $\int_{x1}^{x2}$ .

*Robert Goodhand*

## Basic Definitions

### Mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

or more simply just  $\bar{x} = \frac{1}{n} \sum x_i$

For ungrouped data  $\bar{x} = \frac{\sum f(x)}{\sum f}$

for grouped data we can only estimate the mean by taking the midpoint of each interval to represent that interval midpoint =  $\frac{1}{2}$  (lower class boundary + upper class boundary)

so statistically we do live in a classless society!

### Range

range = highest value – lowest value ignoring correction factors.

### Mean Deviation

mean deviation =  $\sum |x_i - \bar{x}| / n$

and for a frequency distribution

mean deviation =  $\sum |f_i (x_i - \bar{x})| / \sum f_i$

This measure is not widely used.

### Standard Deviation

$$s = \sqrt{\left\{ \sum (x_i - \bar{x})^2 / n \right\}}$$

### Variance

variance = (standard deviation)<sup>2</sup>

$$s^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

$$= \frac{1}{n} \sum (x_i^2 + \bar{x}^2 - 2 x_i \bar{x})$$

$$= \frac{1}{n} (\sum x_i^2 + \sum \bar{x}^2 - 2 \bar{x} \sum x_i)$$

$$= \sum x_i^2 / n + \sum \bar{x}^2 / n - 2 \bar{x} \sum x_i / n$$

$$= \sum x_i^2 / n + \bar{x}^2 - 2 \bar{x}^2$$

Hence we derive the key relationship

$$s^2 = \sum x_i^2 / n - \bar{x}^2$$

### Worked Example

Find the mean and standard deviation of the first n integers

$$\bar{x} = \frac{1}{n} \sum x_i$$

$$= \frac{1}{n} \frac{1}{2} n (n+1)$$

$$= \frac{1}{2} (n+1)$$

This result is used when determining the midpoint term of n terms

$$s^2 = \sum x_i^2 / n - \bar{x}^2$$

$$= \frac{1}{6} n (n+1)(2n+1) / n - \left[ \frac{1}{2}(n+1) \right]^2$$

which reduces to  $\frac{1}{12} (n^2 - 1)$

### Scaling

If each number in a data set is increased by a constant c then

$$\bar{x}_n = \bar{x}_0 + c \text{ but}$$

$$s_n = s_0$$

so adding c affects the mean but leaves the standard deviation unchanged

If each number in the data set is multiplied by constant a

$$\bar{x}_n = a \bar{x}_0 s_0$$

$$s_n = a s_0$$

To summarise if

$$y = ax + b \text{ then}$$

$$\bar{y} = a \bar{x} + b$$

$$s_y = a s_x$$

### Coding

$$\text{if } y = (x - a) / b$$

$$x = a + by$$

$$\bar{x} = a + b \bar{y}$$

$$s_x = b s_y$$

### Interquartile Range (IQR)

$$\text{IQR} = Q_3 - Q_1$$

On a normal distribution this occurs at approximately  $\pm \frac{2}{3} s$

$$Q_1 = \frac{1}{4} (n + 1)^{\text{th}} \text{ value}$$

$$Q_3 = \frac{3}{4} (n + 1)^{\text{th}} \text{ value}$$

### Semi Interquartile Range

$$\text{SIQR} = \frac{1}{2} (Q_3 - Q_1)$$

### Median

The median value is the  $\frac{1}{2} (n+1)^{\text{th}}$  observation.

## Histograms

Frequency density

$$= \text{frequency} / \text{class}$$

width

if the intervals are of equal width

height "bar" = frequency

if the intervals are of unequal width

height "bar" = frequency density

and finally note that area  $\propto$  frequency

### Skew

For positive skew in general

mode < median < mean

For negative skew in general

mode > median > mean

Principle comparators are mode and mean

### Coefficient of Skew

$$\text{mean} - \text{mode} \approx 3 (\text{mean} - \text{median})$$

which leads to

Pearson's coefficient

$$= \frac{3 (\text{mean} - \text{median})}{\text{std. deviation}}$$

$$= \frac{(\text{mean} - \text{mode})}{\text{std. deviation}}$$

The quartile coefficient of skew is given by  $\{(Q_3 - Q_2) - (Q_2 - Q_1)\} / (Q_3 - Q_1)$

### Outliers

Outliers are defined as values that lie

more than  $\frac{3}{2} (Q_3 - Q_1)$

above  $Q_3$  or below  $Q_1$

For a normal distribution this approximates to  $\pm 2 \frac{2}{3} \times \text{std. dev.}$

### **Chebychev's Theorems**

1) The mean and median cannot differ by more than one standard deviation.

This sets the coefficient of skew measure at  $\pm 3$

2) For ANY distribution the proportion of the population that lies outside  $k$  standard deviations is less than  $1/k^2$

This sets the boundary condition for any unknown distribution.

In probabilistic terms we say

$$\Pr (|X - \mu| \geq k\sigma) \leq 1/k^2$$

where  $s$  is the population std. dev.

### **Arrangements**

Number of ways of arranging  $n$  unlike objects in a line is  $n!$

Number of ways of arranging  $n$  objects of which  $p$  are alike is  $n! / p!$

Number of ways of arranging  $n$  objects of which  $p$  are of one type,  $q$  are of another type etc. is  $n! / p! q! \dots$

Number of ways of arranging  $n$  unlike objects in a ring where clockwise/anticlockwise are taken as different arrangements is  $(n - 1)!$

Number of ways of arranging  $n$  unlike objects in a ring where clockwise/anticlockwise are taken as the same arrangement is  $\frac{1}{2} (n - 1)!$

### **Permutations**

Permutations are ordered subsets

$${}^n P_r = n! / (n - r)!$$

### **Combinations**

Combinations are unordered subsets

$${}^n C_r = n! / (n - r)! r!$$

### **Probability Distributions**

we use capital letters to represent the variable  $X$

### **Discrete Random Variables (drv)**

we use lower case letters to represent specific discrete values that the variable  $X$  can take

$$x_1, x_2, x_3, \dots, x_n$$

with probabilities of occurrence as

$$p_1, p_2, p_3, \dots, p_n$$

$X$  is defined as a discrete random variable if and only if

$$p_1 + p_2 + p_3 + \dots + p_n = 1$$

$$\sum p_i = 1 \text{ where } i = 1, 2, 3, \dots, n$$

$$\sum_{\text{all } x} P(X=x) = 1$$

The function that allocates the probabilities is termed the probability density function (pdf).

## Expectation

The expectation of expected value is defined as

$$E(X) = \sum_{\text{all } x} xP(X=x) \text{ or } \sum x_i P_i$$

For symmetrical or uniform distributions  $E(X)$  is the midpoint value.

The concept can be extended to any function of the random variable.

Let  $g(X)$  be any function

$$E[g(X)] = \sum_{\text{all } x} g(x) P(X=x)$$

This relationship has the following colloraries

Cor1:  $E(a) = a$

Cor.2:  $E(aX) = a E(X)$

Cor.3  $E(aX+b) = aE(X) + b$

Cor.4

$$E(f_1(X)+f_2(X)) = E[f_1(x)]+E[f_2x]$$

or simply  $E(X+Y) = E(X) + E(Y)$

$X$  and  $Y$  need not be independent.

## Variance

Taking the experimental approach

For a frequency distribution with mean

$\bar{x}$  and variance  $s^2$  where

$$s^2 = \frac{\sum f(x-\bar{x})^2}{\sum f}$$

which can be written as

$$s^2 = \frac{\sum f x^2}{\sum f} - \bar{x}^2$$

Taking the theoretical approach

For a discrete random variable  $X$  with

$E(X) = \mu$  the variable is defined as

$$V_{ar}(X) = E(X-\mu)^2 = E(X^2) - \mu^2$$

By convention we write this as

$$V_{ar}(X) = E(X^2) - E^2(X) \quad \text{pretty neat}$$

Cor1:  $V_{ar}(a) = 0$

Cor.2:  $V_{ar}(aX) = a^2 V_{ar}(X)$

Cor.3  $V_{ar}(aX + b) = a^2 V_{ar}(X)$

where  $a$  and  $b$  are any constants.

## Cumulative Distribution Function

Let  $X$  be a drv with  $P(X=x)$

then the cdr is given by

$$F(t) = P(X \leq t) = \sum_{x \leq t} P(X=x)$$

$$t = x_1, x_2, x_3, \dots, x_n$$

## Worked example

$X$  is the score from an unbiased die

$$F(t) = t/6 \text{ where } t = 1, 2, 3, 4, 5, 6$$

It is not necessary that there be a specific formula for the cdf.

## Two Random Variable

For  $X$  and  $Y$  independent

$$V_{ar}(X+Y) = V_{ar}(X) + V_{ar}(Y)$$

$$E(aX+bY) = aE(X) + bE(Y) \text{ but}$$

$$V_{ar}(aX \pm bY) = a^2 E(X) \pm b^2 E(Y).$$

If  $X_1$  and  $X_2$  are independent

observations from distribution  $X$  then

$$E(X_1 + X_2) = 2 E(X)$$

$$V_{ar}(X_1 + X_2) = 2 V_{ar}(X)$$

For  $n$  independent observations

$$E(X_1+X_2+X_3\dots X_n) = nE(X) \text{ and}$$

$$\text{Var}(X_1+X_2+X_3\dots X_n) = n\text{Var}(X)$$

In any problem we need to think carefully about whether we are investigating multiples or sums.

Multiples

Sums

$$E(2X) = 2E(X) \quad E(X_1+X_2) = 2E(X)$$

$$\text{Var}(2X) = 4\text{Var}(X) \quad \text{Var}(X_1+X_2) = 2\text{Var}(X)$$

and the relationship n measurements follows on.

### The Binomial Distribution

Where we conduct n independent trials and p is the probability of success in the outcome of any one trial

$$X \sim \text{Bin}(n,p)$$

$$P(X=x) = {}^n C_x p^x q^{n-x} \text{ where } p+q = 1$$

$P(X=x)$  is thus given by the binomial expansion

$$(p+q)^n = {}^n C_0 p^0 q^n + {}^n C_1 p^1 q^{n-1} + \dots + {}^n C_r p^r q^{n-r} + \dots + {}^n C_n p^n q^0$$

Hence the following relationships apply

$$E(X) = np$$

$$\text{Var}(X) = npq$$

$$P(X=r \mid X \sim \text{Bin}(n,p))$$

has the same value as

$$P(X=n-r \mid X \sim \text{Bin}(n,1-p))$$

For cumulative distributions

$$P(X \leq r \mid X \sim \text{Bin}(n,p))$$

has the same value as

$$P(X \geq n-r \mid X \sim \text{Bin}(n,1-p))$$

The mode is the highest value term in the binomial expansion.

### Geometric Distribution

Consider performing a number of independent trials with constant probability of success and q of failure  $(p+q)=1$  we define

$$X \sim \text{Geo}(p)$$

where the drv X has pdf of the form

$$P(X=x) = q^{x-1} p$$

$$\text{Cor.1 } P(X > x) = q^x$$

$$\text{Cor.2 } P[(X > a+b) \mid (X > a)] = P(X > b)$$

If  $X \sim \text{Geo}(p)$  then

$$E(X) = 1/p \text{ and}$$

$$\text{Var}(X) = q/p^2$$

### The Poisson Distribution

If an event is randomly scattered in time or space and has a mean number of occurrences  $\lambda$  in a given interval and X is the random variable

“the number of occurrences”

then  $X \sim \text{Po}(\lambda)$  and

$$P(X=x) = e^{-\lambda} \lambda^x / x!$$

$$x = 0, 1, 2, 3, \dots, \infty \text{ and } \lambda \in \mathbb{R}^+$$

$$\text{we know } e^\lambda = 1 + \lambda + \lambda^2/2! + \dots$$

$$\text{then clearly } \sum P_i = e^{-\lambda} \times e^\lambda = 1$$

so we have a valid pdf.

In general we define

$$E(X) = \sum_{\text{all } x} xP(X=x)$$

so where  $X \sim P_o(\lambda)$

$$\begin{aligned} E(X) &= 0 \times e^{-\lambda} \lambda^0 / 0! + 1 \times e^{-\lambda} \lambda^1 / 1! \\ &\quad + 2 \times e^{-\lambda} \lambda^2 / 2! + 3 \times e^{-\lambda} \lambda^3 / 3! \dots \\ &= \lambda e^{-\lambda} (1 + \lambda + \lambda^2 / 2! + \lambda^3 / 3! \dots) \end{aligned}$$

$$E(X) = \lambda \text{ (ie the mean)}$$

$$\text{Var}(X) = E(X^2) - E^2(X)$$

It can be shown  $E(X^2) = \lambda + \lambda^2$

$$\text{Var}(X) = \lambda + \lambda^2 - \lambda^2 = \lambda$$

(ie both mean and variance are  $\lambda$ )

The Poisson distribution approximates

to the Binomial distribution under

certain conditions

If  $X \sim \text{Bin}(n, p)$

$$P(X=x) = {}^n C_x p^x q^{n-x}$$

$$\text{Now } p = \lambda/n, q = 1 - \lambda/n$$

$$\begin{aligned} \text{so } P(X=x) &= {}^n C_x (\lambda/n)^x (1 - \lambda/n)^{n-x} \\ &= \{ {}^n C_x / n^x \} \lambda^x (1 - \lambda/n)^{n-x} \end{aligned}$$

which we wish to equate to  $1/x! \lambda^x e^{-\lambda}$

For large  $n$   $(1 - \lambda/n)^{n-x} \approx e^{-\lambda}$

and  ${}^n C_x / n^x \approx 1/x!$

and hence the approximation is valid

Examples for  $n = 100$   $p = 0.01$   $\lambda = 1$

x	$P(X=x) \sim \text{Bin}(n, p)$	$P(X=x) \sim P_o(\lambda)$
0	0.3660	0.3697
1	0.3697	0.3697
2	0.1849	0.1839
3	0.0160	0.0613
4	0.0149	0.0153

which is close agreement

Examples for  $n = 50$   $p = 0.1$   $\lambda = 5$

x	$P(X=x) \sim \text{Bin}(n, p)$	$P(X=x) \sim P_o(\lambda)$
	0	0.00515
1	0.0286	0.0337
2	0.0779	0.0842
3	0.139	0.140
4	0.181	0.175

which isn't quite so close

One must also question whether all this is quite so relevant with modern computing power.

### Mode of Poisson Distribution

The mode is the value most likely to occur, that is the one with the highest probability. If  $\lambda$  is an integer the distribution is bimodal with modes at  $x = \lambda - 1$  and  $x = \lambda$

This can readily be seen from Poisson tables.

If  $\lambda$  is not an integer then the mode  $m$  occurs at  $\lambda - 1 < m < \lambda$



## Two Independent Variables

If  $X \sim P_o(m)$

and  $Y \sim P_o(n)$

then  $X + Y \sim P_o(m+n)$

## Probability Distributions

A continuous random variable is a theoretical representation of a continuous variable such as height, mass or time.

A continuous random variable is specified by its probability density function  $f(x)$ .

The pdf is represented by a “curve” and the probabilities are the areas under the curve.

We are therefore concerned with some particular range and we discount any difference between say  $\leq$  and  $<$ .

$$\int_{\text{all } x} f(x) dx = 1$$

$$P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} f(x) dx$$

### Mean

$$E(X) = \int_{\text{all } x} x f(x) dx = \mu.$$

If  $f(x)$  is symmetrical then  $\mu$  will be the  $x$ -value on the line of symmetry.

$$E[g(x)] = \int_{\text{all } x} g(x) f(x) dx$$

$$E[X^2] = \int_{\text{all } x} x^2 f(x) dx$$

Cor.1  $E(a) = a$

Cor.2  $E(aX) = a E(x)$

Cor.3  $E(aX+b) = aE(x) + b$

Cor.4  $E[f_1(X)+f_2(X)]$

$$= E[f_1(x)] + E[f_2(x)]$$

### Variance

If  $X$  is a continuous random variable with probability density function  $f(x)$  then

$$\text{Var}(X) = \int_{\text{all } x} x^2 f(x) dx - \mu^2$$

where  $\mu = E(X) = \int_{\text{all } x} x f(x) dx$

### Standard Deviation

$$\sigma = \sqrt{\text{Var}(X)}$$

Cor.1  $\text{Var}(a) = 0$

Cor.2  $\text{Var}(aX) = a^2 \text{Var}(X)$

Cor.3  $\text{Var}(aX+b) = a^2 \text{Var}(x)$

### Mode

The mode is the value of  $X$  for which  $f(x)$  is a maximum.

### Cumulative Distribution Function

If  $X$  is a continuous random variable with probability density function  $f(x)$  defined for  $a \leq x \leq b$

then the cumulative distribution function is given by

$$F(t) = P(X \leq t) = \int_a^t f(x) dx$$

$a$  and  $b$  may be  $-\infty$  to  $+\infty$

$$P(x_1 \leq X \leq x_2) = F(x_2) - F(x_1)$$

### Median

The median splits the curve into two equal area halves

$$\text{so } \int_a^m f(x) dx = 0.5 \text{ or simply } F(m) = 0.5$$

## Relationship pdf to cdf

$$f(x) = \frac{d}{dx} F(x) = F'(x)$$

ie the gradient of  $F(x)$  gives the value  $f(x)$

## Rectangular Distribution

The distribution has constant height  $1/(b-a)$  between values  $a$  and  $b$  giving an area of 1.

$$\text{So } f(x) = 1/(b-a) \text{ for } a \leq x \leq b$$

$$X \sim R(a, b)$$

$$\text{Hence } \int_{\text{all } x} f(x) dx = \int_a^b 1/(b-a) dx$$

If  $X \sim R(a, b)$  then

$$E(X) = 1/2 (a+b)$$

$$\text{Var}(X) = 1/12 (b-a)^2 \text{ deduced earlier}$$

$$E(X^2) = 1/3 (b^2 + ab + a^2)$$

## Exponential Distribution

A continuous random variable  $X$  having probability density function  $f(x)$  where

$$f(x) = \lambda e^{-\lambda x}$$

is said to be an exponential distribution.

$$\int_{\text{all } x} f(x) dx = \int_0^{\infty} \lambda e^{-\lambda x} dx$$

$$= -[e^{-\lambda x}]_0^{\infty} = e^{-\infty} + e^0 = 1$$

$$P(X < a) = \int_0^a \lambda e^{-\lambda x} dx$$

$$= -[e^{-\lambda x}]_0^a = -e^{-\lambda a} + e^0$$

$$\text{Hence } P(X < a) = 1 - e^{-\lambda a}$$

$$\text{and } P(X > a) = e^{-\lambda a}$$

$$F(x) = P(X \leq x) = 1 - e^{-\lambda x} \text{ for } x \geq 0$$

$$\text{Also } P(X > a+b) | (X > b) = P(X > b)$$

$$E(X) = 1/\lambda$$

$$\text{Var}(X) = 1/\lambda^2$$

$$E(X^2) = 2/\lambda^2$$

The waiting times between successive events in a Poisson distribution follow an exponential distribution.

## Median

Let the median lie at  $x = m$

$$F(m) = 0.5 = 1 - e^{-\lambda m}$$

$$\text{so } e^{\lambda m} = 2 \text{ and hence } m = \ln 2 / \lambda$$

## The Normal Distribution

A continuous variable  $X$  having pdf  $f(x)$  where

$$f(x) = 1/(\sigma\sqrt{2\pi}) e^{-(x-\mu)^2/2\sigma^2}$$

has a normal distribution and we write

$$X \sim N(\mu, \sigma^2)$$

$$E(X) = \mu$$

$$\text{Var}(X) = \sigma^2$$

## Proof

$$E(X) = \int_{\text{all } x} x f(x) dx$$

$$= 1/(\sigma\sqrt{2\pi}) \int_{-\infty}^{+\infty} x e^{-(x-\mu)^2/2\sigma^2} dx$$

$$\text{let } t = (x-\mu)/\sigma$$

$$x = t\sigma + \mu$$

$$dx/dt = \sigma$$

$$E(X) = 1/(\sigma\sqrt{2\pi}) \int_{-\infty}^{+\infty} (\mu + \sigma t) e^{-1/2 t^2} \sigma dt$$

$$= \mu/(\sqrt{2\pi}) \int_{-\infty}^{\infty} e^{-1/2 t^2} dt + \sigma/(\sqrt{2\pi}) \int_{-\infty}^{\infty} t e^{-1/2 t^2} dt$$

$$= \mu + \sigma/2\pi [e^{-1/2 t^2}]_{-\infty}^{\infty}$$

$$\text{Here I assume } 1/(\sqrt{2\pi}) \int_{-\infty}^{\infty} e^{-1/2 t^2} dt = 1$$

$$E(X) = \mu$$

A similar extended proof will show

$$\text{Var}(X) = \sigma^2$$

If  $X \sim N(\mu, \sigma^2)$  the maximum value  $f(x)$  occurs at  $x = \mu$  and  $f(x)$  has points of inflexion at  $x = \mu \pm \sigma$

### Standardised Normal Variable Z

$$Z \sim N(0, 1)$$

Tables give cumulative probabilities – that is areas under the curve.

The symbol for cumulative probability is  $\Phi(z)$  (*pronounced phi*)

$$\text{so } \Phi(z) = P(Z < z)$$

$$\text{if } X \sim N(\mu, \sigma^2)$$

$$\text{and } Z = \frac{X - \mu}{\sigma}$$

$$\text{then } Z \sim N(0, 1)$$

So we have standardised the variable so we only need one universal table.

$$\text{If } Z = \frac{(X - \mu)}{\sigma}$$

$$\text{then } X = \mu + \sigma Z$$

### Approximation Normal to Binomial

if  $X \sim \text{Bin}(n, p)$  with  $n \geq 50$  and  $p \approx 0.5$

then  $X \sim N(np, npq)$

The greater  $n$  the more flexibility on  $p$ .

### Approximation Normal to Poisson

if  $X \sim \text{Po}(\lambda)$

$$\text{then } E(X) = \lambda$$

$$\text{Var}(X) = \lambda$$

for  $\lambda > 20$   $X \sim N(\lambda, \lambda)$

### Random Variables and Sampling

Let  $X$  and  $Y$  be two independent normal variables.

$$X \sim N(\mu_1, \sigma_1^2)$$

$$Y \sim N(\mu_2, \sigma_2^2)$$

then

$$X + Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

$$X - Y \sim N(\mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2)$$

and this may be extended to any number of normal independent variables.

Now take the special case where  $X_1, X_2, X_3 \dots X_n$  are independent observations from same normal distribution

$$X_i \sim N(\mu, \sigma^2) \quad i = 1, 2, 3, \dots, n$$

$$X_1 + X_2 + X_3 + \dots + X_n \sim N(n\mu, n\sigma^2).$$

Remembering if  $X$  is a normal variable

$$X \sim N(\mu, \sigma^2) \text{ then}$$

$$aX \sim N(a\mu, a^2\sigma^2)$$

and it therefore must follow that

$$aX \pm bY \sim N(a\mu_1 \pm b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2)$$

### The Sample Mean

Let  $X_1, X_2, X_3, \dots, X_n$  be a random sample of  $n$  independent observations from a population mean  $\mu$  variance  $\sigma^2$ .

Consider the sample mean  $\bar{X}$

$$\bar{X} = 1/n(X_1 + X_2 + X_3 + \dots X_n)$$

$$= 1/n \sum X_i \quad i = 1, 2, 3, \dots n$$

$$\sum(\bar{X}) = E[1/n(X_1 + X_2 + X_3 + \dots X_n)]$$

$$= 1/n [E(X_1) + E(X_2) \dots E(X_n)]$$

$$= 1/n (n\mu)$$

$$= \mu$$

$$\text{Var}(\bar{X}) = \text{Var}[1/n(X_1 + X_2 + X_3 + \dots X_n)]$$

$$= 1/n^2 [\text{Var}(X_1) + \text{Var}(X_2) \dots \text{Var}(X_n)]$$

$$= 1/n^2 (n\sigma^2)$$

$$\text{Var}(\bar{X}) = \sigma^2/n$$

If we sample without replacement

$$E(\bar{X}) = \mu$$

$$\text{Var}(\bar{X}) = \sigma^2/n \left( \frac{N-n}{N-1} \right)$$

To summarise if we take a random sample of n items from a normal distribution then

$$\bar{X} \sim N(\mu, \sigma^2/n)$$

This is a key result which has even wider implications because under the Central Limit Theorem the distribution can be of any form and still the standard deviation of  $\bar{X}$  is  $\sigma/\sqrt{n}$  and is termed the standard error.

### Distribution of the Sample Proportion

Consider a population in which the proportion of successes is p. If the random sample size is n and X is the

random variable “number of successes”

then  $X \sim \text{Bin}(n, p)$ .

If n is large ( say  $n \geq 50$ )

$$X \sim N(np, npq)$$

Let  $P_s$  be the proportion of successes

$$P_s = X/n$$

$$E(P_s) = E(X/n)$$

$$= 1/n E(X)$$

$$= 1/n (np)$$

$$= p$$

$$\text{Var}(P_s) = \text{Var}(X/n)$$

$$= 1/n^2 \text{Var}(X)$$

$$= 1/n^2 (npq)$$

$$= pq/n$$

Therefore  $P_s \sim N(p, pq/n)$

### Estimation Population Parameters

#### Distribution

#### Parameter

*discrete*

Binomial

n and p

Poisson

$\lambda$

Geometric

P

*continuous*

Rectangular

a and b

Exponential

$\lambda$

Normal

$\mu$  and  $\sigma^2$

If the parameters are unknown we take a sample from the population and use that to make estimates.

This statistic is called an ESTIMATOR

(U,T etc).

The numerical value in any particular instance is called an ESTIMATE (u,t).

Consider a population with unknown parameter  $\theta$

If U is some statistic derived from a random sample taken from the population, then U is an unbiased estimator for  $\theta$  if

$$E(U) = \theta$$

The most efficient estimator is the one which is unbiased and has the smallest variance.

The three important population parameters are mean, variance and proportion (“successes”)

**Parameter Symbol Estimator**

mean	$\mu$	$\mu^{\wedge}$
variance	$s^2$	$\sigma^{\wedge 2}$
proportion	$p$	$p^{\wedge}$

**Parameter Estimator Estimate**

mean	$\bar{X}$	$\bar{x}$
variance	$S^2$	$s^2$
proportion	$P_s$	$p_s$

For mean  $\mu^{\wedge} = \bar{X}$

as  $E(\bar{X}) = \mu$  the estimate is unbiased.

For variance

$$\sigma^{\wedge 2} = \frac{nS^2}{n-1} = \frac{1}{n-1} \sum (X - \bar{X})^2$$

For proportion

$$p^{\wedge} = P_s$$

as  $E(P_s) = p$  and estimate is unbiased.

**Pooled Estimators**

	Size	Mean	Var.
sample 1	$n_1$	$\bar{X}_1$	$S_1^2$
sample 2	$n_2$	$\bar{X}_2$	$S_2^2$

$$\mu^{\wedge} = \{n_1 \bar{X}_1 + n_2 \bar{X}_2\} / \{n_1 + n_2\}$$

This calculation is exact for any set of figures representing sample 1 and sample 2.

$$\sigma^{\wedge 2} = \{n_1 S_1^2 + n_2 S_2^2\} / \{n_1 + n_2 - 2\}$$

This calculation however is dependent on the assumption that the samples are drawn from the same parent population.

	Size	Proportion
sample 1	$n_1$	$P_{s1}$
sample 2	$n_2$	$P_{s2}$

$$p^{\wedge} = \{n_1 P_{s1} + n_2 P_{s2}\} / \{n_1 + n_2\}$$

	one sample	two sample
$\mu^{\wedge}$	$\bar{X}$	$\{n_1 \bar{X}_1 + n_2 \bar{X}_2\} / \{n_1 + n_2\}$
$s^{\wedge 2}$	$n / (n-1) S^2$	$\{n_1 S_1^2 + n_2 S_2^2\} / \{n_1 + n_2 - 2\}$
$p^{\wedge}$	$P_s$	$\{n_1 P_1^2 + n_2 P_2^2\} / \{n_1 + n_2\}$

## Interval Estimation

An interval estimation of an unknown population parameter is an interval constructed so that it has a given probability of including the parameter eg  $P(a < \theta < b) = 0.95$

is termed

the 95% confidence interval for  $\theta$ .

What it is not is the probability that  $\theta$  lies in the interval because it is  $\theta$  that is fixed – albeit unknown – and the interval that varies.

There are three cases to consider

### Case 1

Find the confidence interval for  $\mu$  when the variance  $\sigma^2$  is known.

If  $X \sim N(n, \sigma^2)$

then for any  $n$

$\bar{X} \sim N(\mu, \sigma^2/n)$

If we do not know if  $X$  follows a normal distribution then we need  $n > 30$ . We then standardise

$Z = \{\bar{X} - \mu\} / \{\sigma/\sqrt{n}\}$

so  $Z \sim N(0,1)$

For 95% confidence we set 1.96.

So  $P(\bar{X} - 1.96 \sigma/\sqrt{n} \leq \mu \leq \bar{X} + 1.96 \sigma/\sqrt{n})$

or in brief we write

95% confidence interval is  $\bar{x} \pm 1.96 \sigma/\sqrt{n}$

## Case 2

Find the confidence interval for  $\mu$  when the variance  $\sigma^2$  is unknown.

So first we have to generate an estimate for  $\sigma^2$

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n-1} S^2 \quad (\text{sample variance}) \\ &= \frac{\sum (x - \bar{x})^2}{(n-1)}\end{aligned}$$

95% confidence interval is

$$\bar{x} \pm 1.96 \hat{\sigma} / \sqrt{n}$$

## Case 3

Confidence interval for  $p$  is

$$P_s \pm 1.96 \sqrt{(p_s q_s / n)}$$

so we don't have any prior knowledge of the parent population, the trade off being we require  $n$  large – say  $n > 30$ .

## Significance Testing

The relevant standard deviations for given confidence limits are

	One Tailed	Two Tailed
90%	1.28	1.64
95%	1.64	1.96
99%	2.33	2.58
99.5%	2.58	2.91
1 sd includes	84.1%	68.3%
2 sd includes	97.7%	95.5%
3 sd includes	99.9%	99.7%



