# SAMPLING AND STATISTICAL INFERENCE

## SECTION 4 — FURTHER SINGLE SAMPLE TESTS

### TESTING THE VALUE OF A POPULATION VARIANCE

Suppose we have a random sample of size $n$ from $N(\mu, \sigma^2)$
We want to test

$$H_0 : \sigma^2 = \sigma_0^2 \qquad \text{against} \qquad H_A : \sigma^2 \neq \sigma_0^2$$

The test statistic is $\dfrac{(n-1)s^2}{\sigma_0^2}$

<u>If $H_0$ is true</u>, this test statistic has a <u>known distribution</u> $\left(\text{i.e. } \chi_{n-1}^2\right)$

(For large sample-sizes, the test works quite well even when the population distribution is <u>not</u> Normal.)

### EXAMPLE :

A sample of 25 children (girls, all aged 12) was taken, and their heights measured. The sample mean was 130.3, with sample standard deviation = 3.87. Carry out a test of the hypothesis that $\sigma = 3$ (cm.)
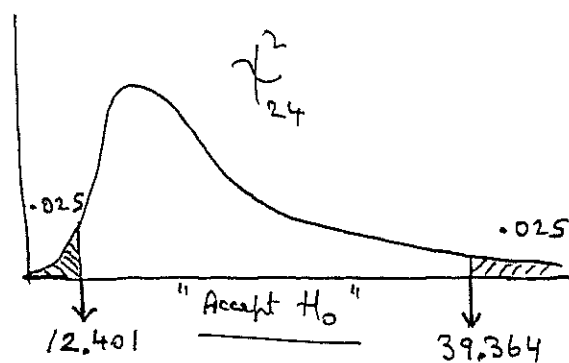
---

$$H_0 : \sigma = 3 \qquad \text{against} \qquad H_A : \sigma \neq 3$$

Now, under $H_0$, $\dfrac{24\,s^2}{3^2}$ is $\chi_{24}^2$ $\qquad (n = 25)$

and $\dfrac{24\,s^2}{3^2} = 39.94$

Since the value of the test statistic falls into the 'CRITICAL REGION', (beyond 39.364) we should <u>reject</u> $H_0$.

(Significance level = 5%)

## TESTING THE VALUE OF THE MEAN OF A Poisson Distribution

Suppose we have a random sample of size $n$ from a Poisson $(\lambda)$ distrib. We want to test

$$H_0 : \lambda = \lambda_0 \quad \text{against} \quad H_A : \lambda \neq \lambda_0$$

The test statistic is the sample sum $\sum X_i$ (where $X_i$ is the Poisson count from $i^{TH}$ sample member)

If $H_0$ is true, this test statistic has a known distribution (i.e. Poisson $(n\lambda_0)$)

NOTE: If $n$ is large (or when $n\lambda_0$ is large) it is possible to use a Normal approximation for the Poisson distribution.

(which is that $\sum X_i \sim$ Poisson $(n\lambda) \longrightarrow N(n\lambda, n\lambda)$ )

If the Normal approximation is used, one would use $\overline{X}$ as the test statistic, with distribution (under $H_0$) as follows:

$$\frac{\overline{X} - \lambda_0}{\sqrt{\lambda_0/n}} \sim N(0,1) \quad \text{, under } H_0 \text{ .}$$

$$\left( \text{or alternatively} \quad \frac{\sum X_i - n\lambda_0}{\sqrt{n\lambda_0}} \sim N(0,1) \right)$$

---

### EXAMPLE :

In an investigation of the frequency of claims by motorists, it was found that there were 960 claims for 6000 policies. Assuming that the number of claims by individual motorists has a Poisson $(\lambda)$ distribution, carry out a test (at the 1% level) of the null hypothesis that $\lambda = 0.17$ against $H_A : \lambda < 0.17$.

Test Statistic $\sum X_i$ is Poisson.

The P-value is Prob $\left[ \sum X_i \leq 960 \mid \text{Poisson mean} = 1020 \right]$

It's necessary to use the Normal approximation.

Since a Poisson variable is discrete, a CONTINUITY CORRECTION is necessary

$$
\begin{aligned}
\text{P-value} &= P\left[ \sum X_i \leq 960.5 \mid (\sum X_i) \text{ is } N(1020, 1020) \right] \\
&= P\left[ \frac{\sum X_i - 1020}{\sqrt{1020}} \leq \frac{960.5 - 1020}{\sqrt{1020}} \right] \\
&= P\left[ Z \leq -1.863 \right] = \underline{0.0312}
\end{aligned}
$$

Thus we could reject $H_0$, if we are prepared to use a significance level of 4% .

# STATISTICAL MODELLING FOR FREQUENCY DISTRIBUTIONS

WE BEGIN WITH AN EXAMPLE:
RUTHERFORD & GEIGER (1910) COLLECTED DATA ON THE NUMBER OF $\alpha$ PARTICLES DETECTED (FROM A RADIOACTIVE SOURCE) FOR EACH OF 2612 INTERVALS OF TIME — EACH OF $7\frac{1}{2}$ SECONDS.

THE FREQUENCY DISTRIBUTION IS AS FOLLOWS

| COUNT | OBSERVED | EXPECTED |
|-------|----------|----------|
| 0 | 57 | 54 |
| 1 | 203 | 210 |
| 2 | 383 | 407 |
| 3 | 525 | 525 |
| 4 | 532 | 509 |
| 5 | 408 | 395 |
| 6 | 273 | 255 |
| 7 | 139 | 141 |
| 8 | 49 | 68 |
| 9 | 27 | 30 |
| 10 | 10 | 11 |
| 11 | 4 | 4 |
| 12 | 2 | 1 |
| $\geq 13$ | 0 | 1 |

IT IS CLEAR THAT THE NUMBER OF DETECTIONS PER INTERVAL VARIES A LOT, AND THERE IS A THEORY WHICH SUGGESTS THAT THE NUMBER OF DETECTIONS HAS A POISSON DISTRIBUTION. FOR WHICH

$$P(X=k) = \frac{m^k e^{-m}}{k!}$$

FOR $k = 0, 1, 2, 3, \ldots$

$X$ IS THE RANDOM VARIABLE (# OF DETECTIONS, HERE) AND $m$ IS A PARAMETER (THE POISSON MEAN).

$m$ CAN BE ESTIMATED BY SIMPLY GETTING THE AVERAGE NUMBER OF DETECTIONS (PER INTERVAL) — AND THIS IS $\hat{m} = 3.877$

THE COLUMN LABELLED EXPECTED IS FOUND BY COMPUTING $\left[\dfrac{(\hat{m})^k e^{-m}}{k!}\right] 2612$ FOR EACH $k = 0, 1, 2, \ldots$

IF THE POISSON MODEL IS CORRECT, THESE ARE THE FREQUENCIES THAT WOULD BE EXPECTED.

WE ASK THE FOLLOWING QUESTION:

Does the Poisson Model fit the observed frequency data?

WE CAN ANSWER QUESTIONS LIKE THIS USING THE FAMOUS CHI-SQUARE GOODNESS-OF-FIT TEST DEVELOPED BY KARL PEARSON.

Notation: ① Denote the observed frequencies by $O_i$

② Denote the expected frequencies by $E_i$
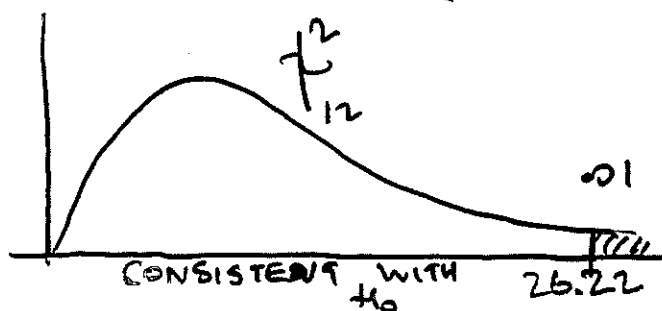
Then we compute the chi-square test statistic

$$\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i}$$

If the $O_i$ and $E_i$ differ greatly, this will be large —— if they are close, this $\chi^2$ will be smaller. For our data, $\chi^2 = 12.4$

Pearson established the sampling distribution of $\chi^2$ 'under $H_0$' (i.e. assuming $H_0$ true —— which in this case means assuming that the observations are from a Poisson distrib.)

This distrib is the chi-square distrib, and its parameter (degrees of freedom) is:

$$D.O.F. = (\#\ \text{Categories}) - 1 - (\#\ \text{Parameters estimated}) = 14-2 = 12$$



CONSISTENT WITH $H_0$    26.22

Conclusion: Freq. data is consistent with a Poisson distrib.

# AN EXAMPLE OF A GOODNESS-OF-FIT TEST

A survey was made of the numbers of boys among families having five children altogether. In 320 families, the number of boys $R$ occurred with the following frequencies:

| Number of boys | 0 | 1 | 2 | 3 | 4 | 5 | Total |
|---|---|---|---|---|---|---|---|
| Observed number of families, $O_r$ | 8 | 40 | 88 | 110 | 56 | 18 | $320 = N$ |
| Expected frequencies, $E_r$ | 10 | 50 | 100 | 100 | 50 | 10 | 320 |

If births are all independent of one another, and the probability $\pi$ of a male birth is the same from one family to another, $R$ should be binomially distributed with parameters $n = 5$ and $\pi$. First let us suppose that $\pi = \frac{1}{2}$. The null hypothesis is now fully specified: '$R$ is binomial with parameters $n = 5$ and $\pi = \frac{1}{2}$'. This gives the set of expected frequencies

$$E_r = N\Pr(r) = N\frac{n!}{r!(n-r)!}\left(\frac{1}{2}\right)^5, \quad r = 0, 1, \ldots, 5.$$

The set of binomial coefficients $n!/r!(n-r)!$ is 1, 5, 10, 10, 5, 1 and $(\frac{1}{2})^5 = \frac{1}{32}$. The values of $E_r$ are thus 10, 50, 100, 100, 50, 10, and $\sum_{r=0}^{5} E_r = 320$. We have put one linear constraint on the $E_r$, the usual one that their total must equal $N$, the total number of observations.

$$X^2 = \sum_{r=0}^{5} \frac{(O_r - E_r)^2}{E_r} = \frac{(8-10)^2}{10} + \frac{(40-50)^2}{50} + \frac{(88-100)^2}{100} + \frac{(110-100)^2}{100}$$

$$+ \frac{(56-50)^2}{50} + \frac{(18-10)^2}{10} = \frac{68}{10} + \frac{136}{50} + \frac{244}{100} = 11.96.$$

This statistic is based on six pairs $(O_r, E_r)$, and the $E_r$ are subject to one linear constraint, so $X^2$ is approximately $\chi^2$ with $6 - 1 = 5$ degrees of freedom. It is significant at the 5% level (the 5% point for $\chi^2_{(5)}$ is 11.07), so at this level we reject the null hypothesis.

We did not give a specific alternative hypothesis, but simply assumed that if the null hypothesis were not true there would be some other set of probabilities, not given by the binomial distribution with $n = 5$ and $\pi = \frac{1}{2}$, that would be more appropriate. Let us consider more carefully what might happen. It is quite possible that births in a family may not be independent events, but that if the first child is a girl the later children are more likely to be girls. In that case, a basic condition for the binomial is violated and we cannot assume either that $\pi$ is constant for all births or that all observations are independent of one another. No simple model can be set up in such a case. However, another alternative is that the binomial conditions do still hold, with $\pi \neq \frac{1}{2}$. This is easy to deal with, and does also appear to explain many sets of data.

If $\pi \neq \frac{1}{2}$, and there is no theoretical reason which gives the exact value of $\pi$, we must estimate $\pi$ from the data. If the data do follow a binomial distribution, the mean, $\bar{r}$, of the observed data will estimate the mean, $n\pi$, of the distribution. Thus $\bar{r}/n$ will estimate $\pi$. We find $\bar{r} = \frac{860}{320} = \frac{43}{16}$. The estimate of $\pi$ is then $\frac{1}{5} \times \frac{43}{16} = 0.5375$; call this $p$. The set of expected frequencies on the null hypothesis that the observations follow a binomial distribution (with $\pi$ not specified in the hypothesis) is therefore

$$N\Pr(r) = N\binom{5}{r}p^r(1-p)^{5-r}$$

$$= N\frac{5!}{r!(5-r)!}(0.5375)^r(0.4625)^{5-r}, \quad r = 0, 1, \ldots, 5.$$

Hence the values of $E_r$ now are 6·8, 39·3, 91·5, 106·3, 61·8, 14·4. The statistic $X^2$ is calculated for the following table of $O_r$ and $E_r$.

| $r$ | 0 | 1 | 2 | 3 | 4 | 5 | Total |
|-----|-----|------|------|-------|------|------|-----------|
| $O_r$ | 8 | 40 | 88 | 110 | 56 | 18 | 320 $= N$ |
| $E_r$ | 6·8 | 39·3 | 91·5 | 106·3 | 61·8 | 14·4 | (320·1) |

$$X^2 = \frac{(8-6\cdot8)^2}{6\cdot8} + \frac{(40-39\cdot3)^2}{39\cdot3} + \frac{(88-91\cdot5)^2}{91\cdot5} + \frac{(110-106\cdot3)^2}{106\cdot3}$$

$$+ \frac{(56-61\cdot8)^2}{61\cdot8} + \frac{(18-14\cdot4)^2}{14\cdot4} = 1\cdot03.$$

The expected values were calculated subject to *two* constraints this time: as usual, $\sum_r E_r = N$, but also this time the mean value of $r$ calculated using the expected frequencies had to equal the mean using the observed frequencies, because this was the equation that we used to estimate $\pi$. This additional constraint is also linear: $\sum r E_r = \sum r O_r$. Thus $X^2$ will be distributed approximately as $\chi^2$ with $6 - 2 = 4$ degrees of freedom.

Since we need one equation for each parameter to be estimated, and each of these equations imposes a constraint on the $E_r$, another way of counting degrees of freedom is as 'number of cells in table *minus* one for total *minus* one for each parameter estimated'. The value 1·03 is certainly not significant as $\chi^2_{(4)}$ and we shall not reject the null hypothesis that the data were binomially distributed. This result suggests that $\pi$ is greater than $\frac{1}{2}$, but that otherwise the binomial conditions are reasonable.

## Contingency tables

Suppose that two characteristics are observed on each of $N$ members of a sample, and that each characteristic is classified into types rather than having an actual measurement recorded. For example, in a human population, we might record colour of hair and colour of eyes for each of $N$ persons. Eyes would be classified 'brown, green/grey, blue' and hair would be classified 'black, brown, fair, ginger'. A summary table would be drawn up, the column headings giving the categories for eye colour and the row headings those for hair colour. Each cell of the table would give the number of people, among the population of $N$, who had a particular eye colour/hair colour combination. Table 18.1 shows a set of results classified in this way.

**Table 18.1** Contingency table for people classified by hair colour and eye colour (observed frequencies)

| Colour of hair | Colour of eyes | | | |
|----------------|-------|-----------|------|-----------|
| | Brown | Green/grey | Blue | Total |
| Black | 50 | 54 | 41 | 145 |
| Brown | 38 | 46 | 48 | 132 |
| Fair | 22 | 30 | 31 | 83 |
| Ginger | 10 | 10 | 20 | 40 |
| Total | 120 | 140 | 140 | 400 $= N$ |

Do the two characteristics, eye colour and hair colour, tend to go together, or is the colour of a person's hair quite independent of the colour of eyes? It seems quite possible, on genetic grounds, that these two characteristics might *not* be independent. We shall now set up a null hypothesis that hair colour and eye colour are independent, and an alternative hypothesis that they are not. We wish to calculate a table of the expected frequencies on the null hypothesis. On this hypothesis, the ratio of the three eye colours, brown: green/grey: blue, should be the same for each one of the hair colours. That is, the ratio should be the same on each row of the table. If the ratio *is* the same on each row, then the best estimate of it is from the totals of the eye colours, namely 120:140:140. This ratio 120:140:140 should then apply to each individual row of the table, so that on each row

there should be a proportion $\frac{120}{400}$ of the row total who have brown eyes, a proportion $\frac{140}{400}$ with green/grey eyes and a proportion $\frac{140}{400}$ with blue eyes. There were 145 people altogether who had black hair, so on the null hypothesis $\frac{120}{400} \times 145 = 43.50$ of these should have brown eyes, $\frac{140}{400} \times 145 = 50.75$ of these should have green/grey eyes, and 50.75 also should have blue eyes (Table 18.2). Similarly, in the second row of the table, there should be $\frac{120}{400} \times 132$ in the first column and $\frac{140}{400} \times 132$ in the second and also in the third column, to account for all the 132 people having brown hair. The third row of the table is dealt with in the same way.

There is a general rule for finding the expected frequencies from the table of observed frequencies, as follows. Let us call the hair-colour totals (145, 132, 83, 40) in the right-hand margin and the eye-colour totals (120, 140, 140) at the foot of the table the *marginal totals*. The total number of observations is $N$ ($= 400$ in the present example). The expected frequency in the cell in row $i$ and column $j$ of the table is equal to $(1/N) \times$ the marginal total of row $i \times$ the marginal total of column $j$. This gives an alternative derivation of Table 18.2.

**Table 18.2** Expected frequency table for people classified by hair colour and eye colour, assuming these two characteristics are independent

| Colour of hair | Colour of eyes | | | |
| --- | --- | --- | --- | --- |
| | Brown | Green/grey | Blue | Total |
| Black | 43·50 | 50·75 | 50·75 | 145·00 |
| Brown | 39·60 | 46·20 | 46·20 | 132·00 |
| Fair | 24·90 | 29·05 | 29·05 | 83·00 |
| Ginger | 12·00 | 14·00 | 14·00 | 40·00 |
| Total | 120·00 | 140·00 | 140·00 | 400·00 |

We now compare Table 18.2 cell by cell with Table 18.1, using

$$X^2 = \sum_{\substack{\text{all} \\ \text{cells}}} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}.$$

(We write $i, j$ as suffices to indicate that summing goes over all rows and all columns of the table – but not, of course, the marginal totals!). The value of $X^2$ is therefore

$$\frac{(50 - 43.50)^2}{43.50} + \frac{(54 - 50.75)^2}{50.75} + \frac{(41 - 50.75)^2}{50.75} + \frac{(38 - 39.60)^2}{39.60} + \frac{(46 - 46.20)^2}{46.20}$$

$$+ \frac{(48 - 46.20)^2}{46.20} + \frac{(22 - 24.90)^2}{24.90} + \frac{(30 - 29.05)^2}{29.05} + \frac{(31 - 29.05)^2}{29.05}$$

$$+ \frac{(10 - 12.00)^2}{12.00} + \frac{(10 - 14.00)^2}{14.00} + \frac{(20 - 14.00)^2}{14.00} = 6.75.$$

In Table 18.2, there are several linear constraints on the calculated frequencies. On the first row, the three expected frequencies must add to 145, the total observed with black hair; this is one constraint, and there is a similar one in the second row and in the third row. But in the fourth row, there is no freedom at all to the expected frequencies; all are constrained by the need for the expected column frequencies to add up to the observed column totals. The full number of constraints applied in the table of expected frequencies is then 1 (first row) + 1 (second row) + 1 (third row) + 3 (last row) = 6. There are 12 cells in the table; hence the degrees of freedom are $12 - 6 = 6$. Thus $X^2$ is approximately $\chi^2_{(6)}$, so its value in this example is not significant. In spite of one or two noticeable discrepancies between an $O_{ij}$ and its corresponding $E_{ij}$, the whole set of observations gives no ground for rejecting the null hypothesis.

The same process can be applied to a table with any number $r$ of rows and $c$ of columns. *The degrees of freedom of the $\chi^2$ variable which will approximate $X^2$ in this general case are* $(r - 1)(c - 1)$. There are $rc$ cells, there is one constraint for each of the first $(r - 1)$ rows, and there are $c$ constraints for the last row; this leaves $rc - (r - 1) - c = (r - 1)(c - 1)$ degrees of freedom.

Two - Sample Problem:

Exa of such a 2-sample problem:

2 Groups of female rats placed on diets with high & low protein content

— gain in wt between the 28th and 84th days of age was measured for each rat. Results in gms

| High P. | | Low P. | |
|---|---|---|---|
| 134 | 107 | 70 | 94 |
| 146 | 83 | 118 | |
| 104 | 113 | 101 | |
| 119 | 129 | 85 | |
| 124 | 97 | 107 | |
| 161 | 123 | 132 | |

Is there evidence of dietary effect — obtain a 95% C. Int for the difference in mean weight gain between the 2 diets.

| Group 1 | Group 2 |
|---|---|
| $\overline{X_1} = 120$ | $\overline{X_2} = 101$ |
| $S_1^2 = 457.45$ | $S_2^2 = 425.33$ |

The analysis of the 2 sample problem in terms of the construction of Conf. Intervals and tests of Hypotheses are easily obtained for the above situation provided we make the following assumptions

(1) Samples 1 and 2 are random samples
— indept of each other

(2) for Sample 1    $X_i^{(1)} \sim N(\mu_1, \sigma^2)$

      "   "   2    $X_i^{(2)} \sim N(\mu_2, \sigma^2)$    } SAME $\sigma^2$

Conf. Int for Diffce $\mu_1 - \mu_2$

With these earlier assumptions

$$\overline{X}_1 \sim N\left(\mu_1, \frac{\sigma^2}{n_1}\right)$$

$$\overline{X}_2 \sim N\left(\mu_2, \frac{\sigma^2}{n_2}\right)$$

and thus

$$\overline{X}_1 - \overline{X}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2}\right)$$

If $\sigma^2$ were known, could proceed as before to

95% Conf Int : $\qquad \overline{X}_1 - \overline{X}_2 \pm 1.96\,\sigma\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$

Since $\sigma^2$ unknown, we are led to follow the same procedure as before:

i.e. estimate $\sigma^2$

Suggest $\qquad \hat{\sigma}^2 = \dfrac{W_1 S_1^2 + W_2 S_2^2}{W_1 + W_2}$

It can be shown that best choice of $W_1, W_2$ is

$$\hat{\sigma}^2 = \frac{(n_1 - 1) S_1^2 + (n_2 - 1) S_2^2}{n_1 + n_2 - 2}$$

Now $\dfrac{(n_1 - 1) S_1^2}{\sigma^2} \sim \chi^2_{n_1 - 1}$ and $\dfrac{(n_2 - 1) S_2^2}{\sigma^2} \sim \chi^2_{n_2 - 1}$

Thus $\dfrac{(n_1 - 1) S_1^2 + (n_2 - 1) S_2^2}{\sigma^2} \sim \chi^2_{n_1 + n_2 - 2}$

and $\dfrac{(\overline{X}_1 - \overline{X}_2) - (\mu_1 - \mu_2)}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim N(0, 1)$

Thus $\dfrac{(\overline{X}_1 - \overline{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\dfrac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{(n_1 + n_2 - 2)\,\sigma^2}}} \Bigg/ \sigma\sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \sim t_{n_1 + n_2 - 2}$

or
$$\frac{(\overline{X}_1 - \overline{X}_2) - (\mu_1 - \mu_2)}{S\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1 + n_2 - 2}$$

where
$$S = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}$$

## 95% Conf Int for $\mu_1 - \mu_2$

$$\overline{X}_1 - \overline{X}_2 \pm t_{n_1+n_2-2, \, .025} \, S\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$= 19 \pm 2.11\sqrt{(446.12)}\sqrt{\frac{1}{12} + \frac{1}{7}}$$

$$= 19 \pm 2.11\sqrt{100.9}$$

$$= 19 \pm 2.11(10.05)$$

$$= 19 \pm 21.2 \qquad \left[\, -2.2 \longrightarrow 40.2 \,\right]$$

# 2 - SAMPLE TEST

May wish to test

$$H_0 : \mu_1 = \mu_2$$
$$H_1 : \mu_1 \neq \mu_2$$
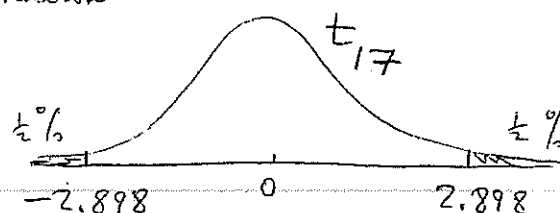
to test for evidence of a dietary effect

Our test statistic will be:

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Assuming $H_0$ true, $t = \dfrac{\bar{X}_1 - \bar{X}_2}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$ and unusually large +ve or -ve values of $t$ would indicate $\mu_1 > \mu_2$ and $\mu_1 < \mu_2$

Distrib of $t$ (under $H_0$) is known

Choose a significance level:
Say 1%

$t_{17}$

$\frac{1}{2}\%$      $\frac{1}{2}\%$

$-2.898$    $0$    $2.898$

Indicate the critical region:

Compute $t = \dfrac{19 - 0}{10.05}$

$$= 1.89$$

Cannot reject $H_0$ at 1% significance level.
    ( nor at the 5% (2.11)
     but Yes at the 10% (1.74)

NOTE

We may accept $H_0$ as being consistent with the data, but it should also be realized that there is a range of values of $(\mu_1 - \mu_2)$ apart from $\mu_1 - \mu_2 = 0$ which are also consistent with the data.

The Extent of this range: Construct the 99% conf int for $(\mu_1 - \mu_2)$
$$\Rightarrow 19 \pm 2.898 (10.05)$$
$$19 \pm 29.1 = -10.1 \text{ to } 48.1$$
$$= -10.1 \text{ to } 48.1$$

Thus any actual difference between the means $(\mu_1 - \mu_2)$ in this range COULD NOT BE DETECTED by our test, using the 1% significance level, i.e. our data is consistent with any actual difference in this range.

Now if some of these differences are considered significant (i.e. important) differences to the experimenter, then this experiment has been completely useless in detecting the presence of such differences.

## TESTING $\sigma_1^2 = \sigma_2^2$

Sample 1 : $X_i^{(1)} \sim N(\mu_1, \sigma_1^2)$    Size $n_1$

" 2 : $X_i^{(2)} \sim N(\mu_2, \sigma_2^2)$    " $n_2$

wish to test $\sigma_1^2 = \sigma_2^2$ (PRELIMINARY STEP IN PREVIOUS ANALYSIS)

We use $\dfrac{(n_1-1)S_1^2}{\sigma_1^2} \sim \chi^2_{n_1-1}$ , $\dfrac{(n_2-1)S_2^2}{\sigma_2^2} \sim \chi^2_{n_2-1}$

Thus $\dfrac{S_1^2}{S_2^2}\dfrac{\sigma_2^2}{\sigma_1^2} \sim F_{n_1-1, n_2-1}$

Under $H_0 : \sigma_1^2 = \sigma_2^2$, we see that

$$\dfrac{S_1^2}{S_2^2} \sim F_{n_1-1, n_2-1}$$

Thus, we can test $H_0$, using the Test Statistic $\dfrac{S_1^2}{S_2^2}$, which will be close to 1 when $H_0$ true.

## TEST PROCEDURE:

Form $\dfrac{S_1^2}{S_2^2}$    so that ratio $> 1$