Semester 2, 2012 (Last adjustments: October 23, 2012)

**Lecture Notes**

# MATH1905 Statistics (Advanced)

**Lecturer**

Dr. John T. Ormerod

School of Mathematics & Statistics F07

University of Sydney

(w) 02 9351 5883

(e) john.ormerod (at) sydney.edu.au

Semester 1, 2012 (Last adjustments: October 23, 2012)

# Lecture 1 - Content

☐ $\chi^2$ **goodness of fit tests**

☐ **Further applications of** $\chi^2$ **GoF tests**

# $\chi^2$ Goodness of fit tests

## Motivational setting

☐ Suppose we have $n$ independent trials with $X$ successes and $n - X$ failures: $X \sim \mathcal{B}(n, p)$.

☐ Test $H_0 : p = p_0$ against $H_1 : p \neq p_0$. If $H_0$ is true then

|  | Success | Failure | Total |
|---|---|---|---|
| Observed | $O_1 = X$ | $O_2 = n - X$ | n |
| Expected | $E_1 = np_0$ | $E_2 = n(1 - p_0)$ | n |

☐ Large values of $|X - np_0|$ support $H_1$. Thus large values of

$$\tau = \frac{(X - np_0)^2}{np_0(1 - p_0)}$$

support $H_1$.

## Motivational setting (cont)

☐ Note,

$$(O_2 - E_2)^2 = [(n - X) - (n - np_0)]^2$$
$$= (X - np_0)^2 = (O_1 - E_1)^2.$$

☐ Also,

$$\frac{1}{np_0} + \frac{1}{n(1 - p_0)} = \frac{1}{np_0(1 - p_0)}.$$

☐ Thus,

$$\tau = \frac{(X - np_0)^2}{np_0(1 - p_0)} = (X - np_0)^2 [\frac{1}{np_0} + \frac{1}{n(1 - p_0)}]$$
$$= \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2}.$$

☐ This is a special case of Pearson's $\chi^2$ statistic.

# Pearson's $X^2$ GoF test

Assume we have $g$ categories, not just success/failure and $H_0$ specifies a model giving expected frequencies for each category.

**Definition 1.** Pearson's $\chi^2$ test-statistic is

$$X^2 = \sum_{i=1}^{g} \frac{(O_i - E_i)^2}{E_i},$$

where $O_i$ is the observed frequency in the $i$th category and $E_i$ is the expected frequency if $H_0$ is true.

## Alternative Calculation Formula for $X^2$

Since $\sum_{i=1}^{g} O_i = n$ and $E_i = np_i$ it follows that:

$$\sum_{i=1}^{g} E_i = \sum_{i=1}^{g} np_i = n$$

Hence,

$$\sum_{i=1}^{g} \frac{(O_i - E_i)^2}{E_i} = \sum_{i=1}^{g} \frac{O_i^2 - 2O_iE_i + E_i^2}{E_i}$$

$$= \sum_{i=1}^{g} \frac{O_i^2}{E_i} - 2O_i + E_i$$

$$= \left( \sum_{i=1}^{g} \frac{O_i^2}{E_i} \right) - 2n + n$$

$$= \sum_{i=1}^{g} \frac{O_i^2}{E_i} - n.$$

# Pearson's $X^2$ GoF test (cont)

☐ Thus, the easiest form for calculation purposes of $X^2$ is,

$$X^2 = \sum_{i=1}^{g} \frac{O_i^2}{E_i} - n.$$

☐ We reject the model if the $X^2$-statistic is too large.

☐ The sampling distribution of the statistic has (asymptotically) a chi-squared distribution with $g-1$ degrees of freedom.

$$P\text{-value} = \mathrm{P}(\chi_{g-1}^2 \geq \text{observed } X^2 \text{ value}).$$

   ○ Note that $\chi_1^2 = Z^2$, where $Z \sim \mathcal{N}(0,1)$. In R with `pchisq()` or tables.

   ○ The $X^2$ test should only be used when the expected frequencies, $E_i$, are greater than 5.
(Recall this corresponds to $np \geq 5$ for the normal approximation to the binomial!)

# The $\chi^2$ distribution

☐ The $\chi^2$ r.v. can only take non-negative values. The distribution is not symmetric but is right skewed. Tables typically give

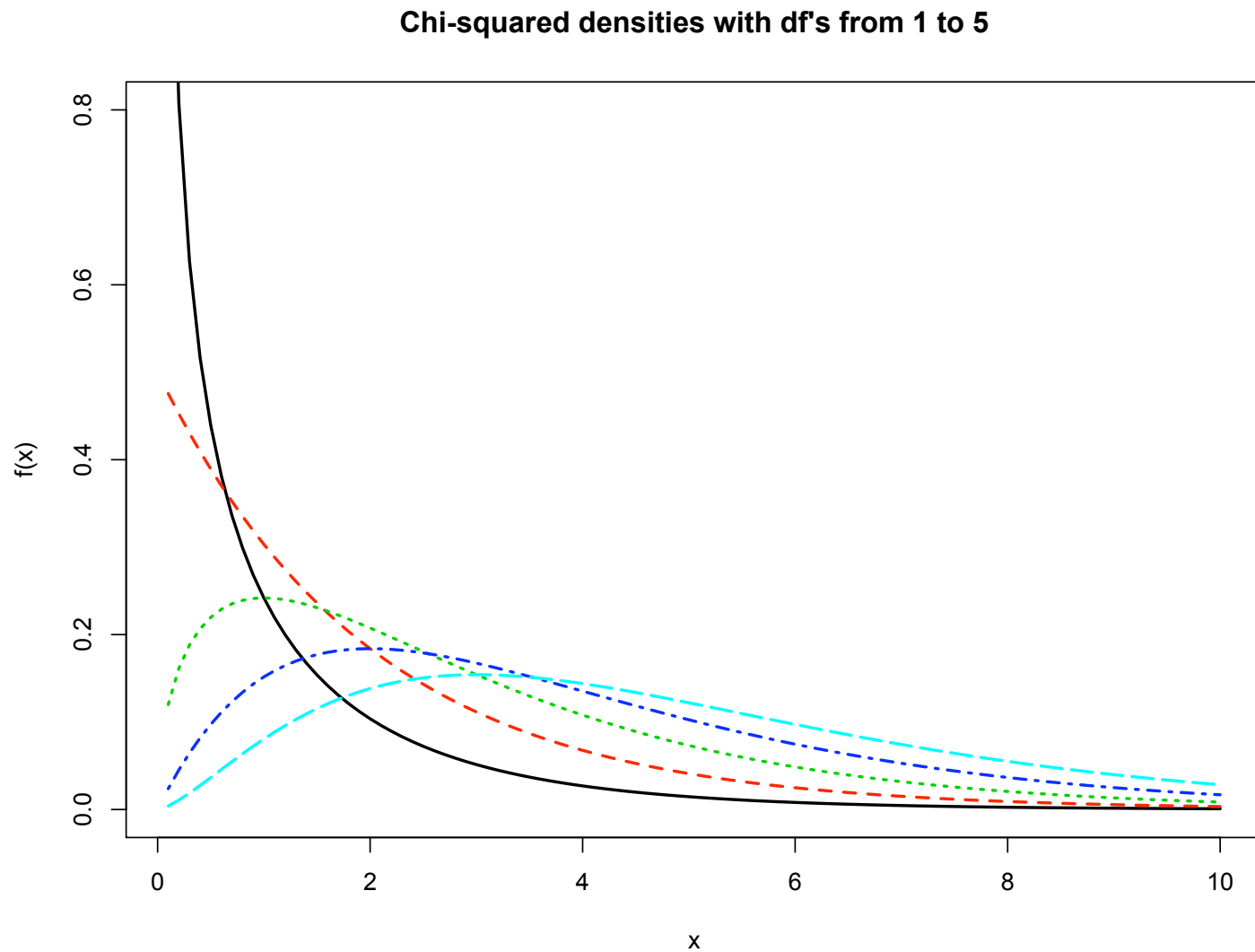$$P(\chi^2_\nu > x) = p = \texttt{1-pchisq(x,nu)}$$

for particular d.f. $\nu$ and $p$ values.

☐ With R we can visualise the densities:

```
> x = 1:100/10
> plot(x,dchisq(x,1),type="l",ylim=c(0,0.8))
> points(x,dchisq(x,2),type="l",lty=2,col=2)
> points(x,dchisq(x,3),type="l",lty=3,col=3)
> points(x,dchisq(x,4),type="l",lty=4,col=4)
> points(x,dchisq(x,5),type="l",lty=5,col=5)
```

# The $\chi^2$ distribution (cont)

**Chi-squared densities with df's from 1 to 5**

## Example.

(i) $P(\chi_1^2 > 3.841) = ?$

  □ From tables the points tabulated below are $x$, where $P(\chi_\nu^2 > x) = p$

| | | | | $p$ | | |
|---|---|---|---|---|---|---|
| $\nu$ | 0.25 | 0.15 | 0.10 | 0.05 | 0.025 | 0.01 |
| 1 | 1.323 | 2.072 | 2.706 | 3.841 | 5.024 | 6.635 |

  □ Using R:

```
> 1-pchisq(3.841,1)
[1] 0.05001368
> pchisq(3.841,1,lower.tail=FALSE)
[1] 0.05001368
```

(Note also that $1.96^2 = 3.841$, $P(|Z|^2 > 1.96^2) = 0.05$.)

(ii) $P(\chi_{10}^2 > 20) = ?$

□ From tables

| $\nu$ | 0.25 | 0.15 | 0.10 | 0.05 | 0.025 | 0.01 |
|---|---|---|---|---|---|---|
| | | | | $p$ | | |
| 10 | 12.549 | 14.534 | 15.987 | 18.307 | 20.483 | 23.209 |

So that $0.025 < P(\chi_{10}^2 > 20) < 0.05$.

□ Via R

```
> pchisq(20,10,lower.tail=FALSE)
[1] 0.02925269
```

(iii) $P(13.848 < \chi_{24}^2 < 39.364) = ?$

□ Can't be obtained because tables only have $1 \le \nu \le 10$.

□ Via R

```
> pchisq(39.364,24) - pchisq(13.848,24)
[1] 0.9250087
```

**Example** (Phenotypes, PQ p117)**.** In an experiment involving a dihybrid cross of flies, 148 progeny were classified by phenotype as follows.

|    | AB | Ab | aB | ab | Total |
|----|----|----|----|----|-------|
|    | 87 | 31 | 25 | 5  | 148   |

☐ Genetic theory predicts a ratio 9:3:3:1 for AB:Ab:aB:ab.

☐ Do the data support the theory?

☐ The model specifies the expected frequencies

| AB | Ab | aB | ab |
|----|----|----|----|
| $\frac{9}{16} \times 148$ | $\frac{3}{16} \times 148$ | $\frac{3}{16} \times 148$ | $\frac{1}{16} \times 148$ |

| $E_i :$ | 83.25 | 27.75 | 27.75 | 9.25 |
|---------|-------|-------|-------|------|

**Example** (cont).

☐ The test statistic is

$$X^2 = \frac{87^2}{83.25} + \frac{31^2}{27.75} + \frac{25^2}{27.75} + \frac{5^2}{9.25} - 148$$

$$= 2.775.$$

☐ The $P$-value for testing fit of the model is

$$P\text{-value} = P(\chi_3^2 \geq 2.775) = \texttt{pchisq(2.775,3,lower.tail=FALSE)}.$$

☐ Since the $P$-value is large, 0.43 (2dp), we conclude that the data are consistent with the model.

**Example** (Accidents). The number of fatal accidents on NSW roads in the July month between 2004 - 2010 were:

| 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 |
|------|------|------|------|------|------|------|
| 44   | 50   | 34   | 41   | 34   | 27   | 29   |

Test the claim that the accident rate didn't change over this seven year period:

☐ $p_i$ denotes the probability that a fatal accident is 'allocated' to month $i$.

☐ Model: $p_i = \frac{1}{7}, \quad i = 1, \ldots, 7$.

☐ The total number of accidents is 259. Thus $E_i = \frac{259}{7} = 37$.

$$\text{The test statistic is } X^2 = \sum_{i=1}^{7} \frac{O_i^2}{E_i} - 259 = 11.24.$$

☐ Thus, the $P$-value $= 0.081 \Rightarrow$ no rejection of the claim that the accident rate is constant across the years based on these data.

# Further applications of $\chi^2$ GoF tests

□ Recall that for known parameters,

$$X^2 = \sum_{i=1}^{g} \frac{(O_i - E_i)^2}{E_i} = \sum_{i=1}^{g} \frac{O_i^2}{E_i} - n \overset{\text{under } H_0}{\sim} \chi^2_{g-1}.$$

□ If we want to check the fit of a model that involves unknown parameters we first have to estimate the parameters.

□ Since we use the same data to estimate the parameters and test the fit we find the sampling distribution of the $X^2$ statistic has to be adjusted.

□ The distribution is still $\chi^2$ but the dfs are reduced to $g - k - 1$, where

　○ $g$ is the number of categories and

　○ $k$ is the smallest number of parameters that need to be estimated using the data.

**Example** (More on phenotypes). In a backcross experiment to investigate the genetic linkage between two factors A and B in a species of flower, some researchers classified 400 offspring by phenotype as follows:

$$\begin{array}{cccc} \text{AB} & \text{Ab} & \text{aB} & \text{ab} \\ 128 & 86 & 74 & 112 \end{array}$$

(i) Under the no linkage model, the four phenotypes are equally likely.
   *Show that this model is a poor fit.*

(ii) If linkage is in the coupling phase, the probabilities of the four phenotypes are

$$\begin{array}{cccc} \text{AB} & \text{Ab} & \text{aB} & \text{ab} \\ \frac{1}{2}(1-p) & \frac{1}{2}p & \frac{1}{2}p & \frac{1}{2}(1-p) \end{array}$$

   where $p$ is the 'recombination fraction' and is estimated by the overall proportion of Ab and aB.
   *Show that this model fits the data well.*

**Example** (cont; (i))**.** Model says that all categories are equally likely.

☐ The total number of observations is $n = 400$.

☐
$$\text{Observed } O_i : \quad 128 \quad 86 \quad 74 \quad 112$$
$$\text{Expected } E_i : \quad 100 \quad 100 \quad 100 \quad 100$$

☐ $X^2 = \sum_i \frac{O_i^2}{E_i} - n = 18$.

☐ $P = \mathrm{P}(\chi_3^2 \geq 18) < 0.01$.

☐ Thus the data are not consistent with the model.

**Example** (cont; (ii)). Here, $\widehat{p} = (86 + 74)/400 = 0.4$.

☐ The expected frequencies are

$$E_1 = E_4 = 400 \times 0.3 = 120$$

$$E_2 = E_3 = 400 \times 0.2 = 80.$$

☐ $X^2 = 1.97$ and $g - k - 1 = 2$

☐ $P = \mathrm{P}(\chi^2_2 \geq 1.97) > 0.10$.

☐ The data are consistent with this model.

Tuesday, 23 October 2012

# Lecture 2 - Content

☐ **Further applications of $\chi^2$ GoF tests (cont)**

**Example** (Infections). 200 groups of 5 insects each were inspected. For each group the number of infected insects $(x)$ was counted giving:

$$3, 2, 5, 1, 0, .., 2.$$

The data were condensed into the table below, writing $x_i$ for the number infected and $f_i$ for the corresponding frequency:

| $x_i$ : | 0 | 1 | 2 | 3 | 4 | 5 | Total |
|---|---|---|---|---|---|---|---|
| $f_i$ : | 20 | 62 | 55 | 38 | 20 | 5 | 200 |

Does the binomial model fit the data?

☐ The null hypothesis is that $X \sim \mathcal{B}(5, p)$. We need to estimate $p$.

☐ There were $5 \times 200 = 1000$ insects in total and 391 of these were infected, i.e. an estimate is

$$\widehat{p} = \frac{391}{1000} = 0.391.$$

## Example (cont).

| $i$ | 0 | 1 | 2 | 3 | 4 | 5 | Total |
|-----|-----|-----|-----|-----|-----|-----|-----|
| $p_i$ | 0.0837 | 0.2689 | 0.3453 | 0.2217 | 0.0712 | 0.0091 | 1.00 |
| $E_i$ | 16.754 | 53.783 | 69.061 | 44.340 | 14.234 | 1.828 | 200 |
| $O_i$ | 20 | 62 | 55 | 38 | 20 | 5 | 200 |

Notice that the $E_i$ value falls below 5 for the last group. We combine the last two cells together.

| $i$ | 0 | 1 | 2 | 3 | $\geq 4$ | Total |
|-----|-----|-----|-----|-----|-----|-----|
| $E_i$ | 16.754 | 53.783 | 69.061 | 44.340 | 16.062 | 200 |
| $O_i$ | 20 | 62 | 55 | 38 | 25 | 200 |

Hence, $X^2 = 10.63$ and because $g = 5$ and $k = 1$ we obtain

$$P = \mathrm{P}(\chi_3^2 \geq 10.63) < 0.025.$$

Thus, the binomial model does not fit.

# Testing the fit of a normal model

Given a data set $x_1, x_2, .., x_n$ we want to test if the data come from a $\mathcal{N}(\mu, \sigma^2)$ population.

(a) First calculate the sample mean, $\overline{x}$, and the sample variance, $s^2$.

(b) Form a grouped frequency table summary of the data with (ideally) 5 to 10 categories. Aim to have at least 5 values in each category.

(c) To check the normal claim work out the expected frequencies for each category by fitting $\mathcal{N}(\overline{x}, s^2)$.

(d) Use $X^2$ as the test statistic. To calculate the $P$-value use $(g - 2 - 1)$ df.

**Example** (Rainfall). We have $n = 30$ observations corresponding to Sydney's annual rainfall (in inches) from 1980-2009 (from http://www.bom.gov.au):

```
> y = c( 956,1083,1499, 994, 816, 995,1200, 860,1359, 822,
+        1470,1649,1078,1149,1230, 907, 913,1282,1121,1977,
+        1526,1862,1313,1225,1217,1801,1346, 838,1038, 736)
> x = 2009:1980
```

Test if the rainfall follows a normal distribution.

(a) $\overline{y} = 1208.733$ and $s^2 = 105762.8$.

(b) Grouping the data into a frequency table:

| Interval | Frequency |
|:---:|:---:|
| $y \leq 900$ | 5 |
| $900 < y \leq 1200$ | 11 |
| $1200 < y \leq 1500$ | 9 |
| $y > 1500$ | 5 |

**Example** (cont).

(c) We now calculate the expected frequencies using

$$Y \sim \mathcal{N}(1208.733, 325.212^2)$$

$$P(Y \leq 900) = \mathrm{P}\left(Z \leq \frac{900 - 1208.733}{325.212}\right) = \mathrm{P}(Z \leq -0.95) = 0.1712.$$

Thus $E_1 = 30 \times 0.1712 = 5.137$.

$$P(900 < Y \leq 1200) = P(-0.95 < Z \leq -0.03)$$
$$= 0.318$$

$$E_2 = 30 \times 0.318 = 9.542.$$

Similarly,

$E_3 = 30 \times 0.3254 = 9.765$ and

$E_4 = 30 - 5.137 - 9.542 - 9.765 = 5.556.$

**Example** (cont). (d) Our table is

| Interval | Frequency | Expected |
|:---:|:---:|:---:|
| $y \leq 900$ | 5 | 5.137 |
| $900 < y \leq 1200$ | 11 | 9.542 |
| $1200 < y \leq 1500$ | 9 | 9.765 |
| $y > 1500$ | 5 | 5.556 |

The test statistic is

$$X^2 = \frac{5^2}{5.137} + \frac{11^2}{9.542} + \frac{9^2}{9.765} + \frac{5^2}{5.556} - 30 = 0.342.$$

Here $g = 4$ and $k = 2$ so we have 1 d.f.

The $P$-value is $P(\chi_1^2 \geq 0.342) = 0.559$, with R.

Thus the data are consistent with the normal model.

# Tests for independence

If we have data classified according to two attributes then we can construct a contingency table which is just a convenient way of presenting the group frequencies. For example, we have data on 422 drivers and motorcyclists killed in NSW in 1988. We classify the people by blood alcohol level and gender.

| Alc (g/100ml) | 0 | $(0, 0.08)$ | $[0.08, 0.15)$ | $\geq 0.15$ | Total |
|---|---|---|---|---|---|
| Male | 206 | 37 | 35 | 76 | 354 |
| Female | 53 | 5 | 4 | 6 | 68 |
| Total | 259 | 42 | 39 | 82 | |

Test the claim that gender is independent of blood alcohol level.

# A probability model for contingency tables

□ Let $p_{ij}$ denote the probability of a victim being gender $i$ and alcohol level group $j$ then the independence model says:

$$p_{ij} = p_i^g p_j^a, \quad \text{where}$$

○ $p_i^g$ is the prob. of being of gender $i$,

○ $p_j^a$ is the prob. of being in alcohol group $j$.

□ We estimate $p_i^s$ and $p_j^a$ by the marginal proportions,

$$\widehat{p}_1^g = \frac{354}{422} \; \Rightarrow \; \widehat{p}_2^g = 1 - \widehat{p}_1^g = \frac{68}{422};$$

$$\widehat{p}_1^a = \frac{259}{422}, \; \widehat{p}_2^a = \frac{42}{422}, \; \widehat{p}_3^a = \frac{39}{422}, \; \text{and} \; \widehat{p}_4^a = \frac{82}{422}.$$

□ The expected frequency under the independence model in the Male/Alcohol 0 group is

$$422 \times \widehat{p}_{11} = 422 \times \frac{259}{422} \times \frac{354}{422} = 217.265.$$

□ For a general table with entries $x_{ij}$ the expected frequencies under independence are

$$E_{ij} = n \times \frac{x_{i\bullet}}{n} \times \frac{x_{\bullet j}}{n} = \frac{x_{i\bullet} \times x_{\bullet j}}{n},$$

where $x_{i\bullet}$ denotes the sum of the $i$th row and $x_{\bullet j}$ denotes the sum of the $j$th column.

□ The expected frequencies for the accident data are

```
217.265    35.232    32.716    68.787
 41.735     6.768     6.284    13.213
```

□ Thus, $X^2 = \sum \frac{O_i^2}{E_i} - n = \ldots = 9.859.$

□ We have $g = 8$ groups in this example and $k = 1 + 3 = 4$ estimated parameters. Thus the df for $X^2$ is $8 - 4 - 1 = 3$. Hence, `1-pchisq(9.859,3)` equals 0.02 (2dp).

□ Thus, there is strong evidence to suggest blood alcohol level and gender are related in accident victims.

□ In general if we have a table with $r$ rows and $c$ columns then the test statistic for testing independence will have

$$(r-1)(c-1) \text{ degrees of freedom.}$$

# Tests for Symmetry in Tables

The following data for 205 married people were reported in Yule (1900).

|  Husband | Wife Tall | Medium | Short |
|---|---|---|---|
| Tall | 18 | 28 | 14 |
| Medium | 20 | 51 | 28 |
| Short | 12 | 25 | 9 |

Here there are $g = 9$ groups. Suppose we wish to test for table symmetry, i.e. the probability of Tall Men marrying Short Women and Short Men marrying Tall Women are roughly equal.

The table corresponding to a symmetric table model would have $k = 5$ parameters

|  | Wife | | |
| Husband | Tall | Medium | Short |
| --- | --- | --- | --- |
| Tall | $p_1$ | $p_2$ | $p_3$ |
| Medium | $p_2$ | $p_4$ | $p_5$ |
| Short | $p_3$ | $p_5$ | $1 - p_1 - \cdots - p_5$ |

# Tests for Symmetry in Tables

| | Wife | | |
|:---:|:---:|:---:|:---:|
| Husband | Tall | Medium | Short |
| Tall | 18 | 28 | 14 |
| Medium | 20 | 51 | 28 |
| Short | 12 | 25 | 9 |

Under the symmetric model the expected values of the table entries would be

| | Wife | | |
|:---:|:---:|:---:|:---:|
| Husband | Tall | Medium | Short |
| Tall | 18 | 24 | 13 |
| Medium | 24 | 51 | 26.5 |
| Short | 13 | 26.5 | 9 |

# Tests for Symmetry in Tables

The value of the test statistic would be

$$X^2 = \frac{(18-18)^2}{18} + \frac{(28-24)^2}{24} + \cdots + \frac{(9-9)^2}{9} = 1.656991$$

The degrees of freedom would be $\nu = g - k - 1 = 9 - 5 - 1 = 3$.

The $P$-value would then be

$$P(\chi_3^2 > 1.656991) = \texttt{1-pchisq(1.656991,3)} = 0.6465376$$

So we would accept the null hypothesis for table symmetry.

Note that for a test for independence the test statistic is $2.907$ (check yourself!) with $\nu = 4$ degrees of freedom and a $P$-value of $0.5735$. Thus, a null hypothesis of independent rows and columns is also consistent with the data!

# THE END !!!

Monday, 24th October 2011

# Lecture 3 - Content

☐ **Extended Answer Section of 2010 Exam**

☐ **Extended Answer Section of 2011 Exam**

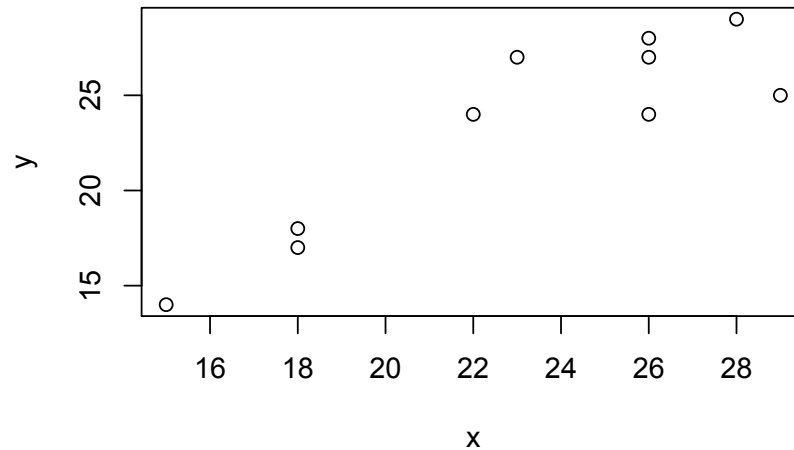## Question 1 (16 marks in all)

The following R-output gives the maximum temperatures $x_i$ and $y_i$, in °C, on two successive days (day 1 and day 2) in 10 Australian weather observation stations, $i = 1, \ldots, 10$:

```
> x = c(22,18,26,26,29,15,23,18,28,26)
> y = c(24,17,28,24,25,14,27,18,29,27)
```

Additionally you might find the following R output of use:

```
> length(x)        > sum(x^2)
[1] 10             [1] 5539
> sum(x)           > sum(y^2)
[1] 231            [1] 5669
> sum(y)           > sum(x*y)
[1] 233            [1] 5580
```



(a) (3 marks) *Calculate the coefficient of correlation.*

$$r_{xy} = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \frac{5580 - 10 \times 23.1 \times 23.3}{\sqrt{(5539 - 231^2/10) \times (5669 - 233^2/10)}}$$

$$= \frac{197.7}{\sqrt{202.9 \times 240.1}} = 0.8957.$$

(b) (4 marks) *Find the equation of the least squares line of temperature on day 2 $(y)$ on temperature of day 1 $(x)$.*

$$b = \frac{S_{xy}}{S_{xx}} = \frac{197.7}{202.9} = 0.9744$$

and

$$a = \bar{y} - b\bar{x} = 23.3 - b \times 23.1 = 0.7920.$$

Thus, the regression line is $y = 0.7920 + 0.9744 \times x$. [2 marks for writing down the regression line explicitly, 1 mark if $a$ correct, 1 mark if $b$ correct].

(c) (2 marks) *The simplest weather forecast is 'tomorrow will be like today'. Use the regression line to give an improved forecast in one of those weather stations, if today's temperature there is $15°C$*

Let $x = 15$, then $y = 0.7920 + 0.9744 \times 15 = 15.408$ degrees.
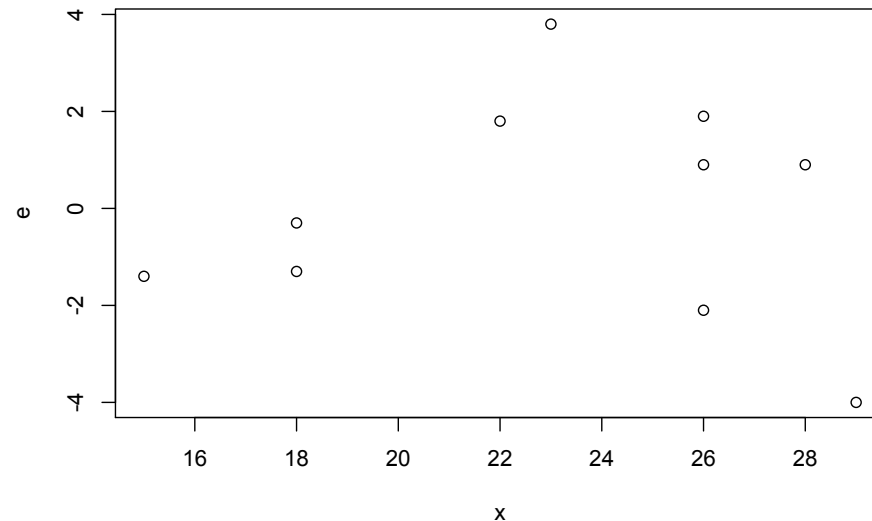
(d) (6 marks) *The residuals $e_1, \ldots, e_{10}$ from the least squares fit are as follows (to 1dp):*

```
> round(lm(y~x)$res,1)
    1     2     3     4     5     6     7     8     9    10
  1.8  -1.3   1.9  -2.1  -4.0  -1.4   3.8  -0.3   0.9   0.9
```

*Draw a residual plot and hence comment on the appropriateness of the linear regression.*

The residual plot should look as follows (deduct marks if axes are not labeled and if the $e$ is not drawn on the $y$-axis etc) [3 marks]



The sample size is only 10 and therefore it is difficult to justify with certainty the appropriateness of the linear regression [1 mark]. There seems to be some 'quadratic' curvature left [1 mark], also the variability of the residuals seem to be larger for larger $x$ values [1 mark].

(e) (1 mark) *About what percentage of the variability of $y$'s is explained by the regression line?*

About 80%, since $r^2 = r_{xy}^2 = 0.8957^2 = 0.8022$. [1 mark]

# Question 2 (15 marks in all)

The following data are measurement of weight gain (in gm) after 10 male rats and 10 female rats were given the same diet over the same period of time. The 10 male rats and 10 female rats were chosen independently.

$$\text{Male } (x) \quad 2.6 \;\; 4.8 \;\; 12.5 \;\; 8.7 \;\; 9.7 \;\; 8.2 \;\; 9.4 \;\; 8.7 \;\; 9.2 \;\; 10.0$$
$$\text{Female } (y) \quad 8.1 \;\; 7.6 \;\; 10.5 \;\; 8.9 \;\; 11.2 \;\; 6.9 \;\; 11.7 \;\; 12.6 \;\; 10.3 \;\; 7.1$$

(a) (8 marks) *Assume that both samples can be modelled by a normal distribution with the same population variance. Given that $\sum_{i=1}^{10} x_i = 83.8$, $\sum_{i=1}^{10} y_i = 94.9$, $\sum_{i=1}^{10} x_i^2 = 771.76$, $\sum_{i=1}^{10} y_i^2 = 938.03$, is there evidence of a difference in weight gains between male rats and female rats?*

Here, a two-sample $t$-test is the appropriate test to be used because male rats and female rats were chosen independently [1 mark].

Our null and alternative hypotheses are [1 mark]:

$$H_0 \colon \mu_x = \mu_y \quad \text{vs} \quad H_1 \colon \mu_x \neq \mu_y$$

Thus, the pooled variance is [1 mark]

$$s_p^2 = \frac{(n_x - 1)s_x^2 + (n_y - 1) \times s_y^2}{18} = \frac{1}{2}(s_x^2 + s_y^2)$$
$$= (7.724 + 4.158778)/2 = 5.941389$$

$$t_{(n_x+n_y-2)} \sim \frac{\overline{X} - \overline{Y}}{s_p\sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} = \frac{8.38 - 9.49}{\sqrt{5.941389}\sqrt{1/5}} = -1.018272, \quad [2 \text{ marks}]$$

The degrees of freedom is $\nu = 18$ [1 mark].

From the tables we find the critical value at the $5\%$ level to be between $-2.228$ and $-2.086$. Alternatively the $P$-value is $P = 2 \times P(t_{18} > 1.018272)$. From tables $P(t_{10} > 0.700) = 0.25$ and $P(t_{20} > 1.064) = 0.15$ so $0.3 < P < 0.5$ [1 mark].

Thus we accept the null hypothesis that there is no difference in weight gains between male rats and female rats. [1 mark]

(b) (7 marks) Let $X \sim \mathcal{N}(\mu, \sigma^2)$ with $\mu = 1$ and $\sigma^2 = 1/2$. *Use Chebyshev's inequality to bound* $\mathrm{P}(X > 2)$.

Chebyshev's inequality is

$$\mathrm{P}(|X - \mu| \geq c\sigma) \leq \frac{1}{c^2} \quad \text{[1 mark]}.$$

Here, $\mu = 1$ and $\sigma = \frac{1}{\sqrt{2}}$ [1 mark].

Thus,

$$\mathrm{P}(X > 2) \leq \mathrm{P}(|X - 1| > 1) \leq \mathrm{P}(|X - 1| > \sqrt{2}\frac{1}{\sqrt{2}}) < 1/2. \quad \text{[3 marks]}.$$

*Is Chebyshev's bound obtained in (i) sharp in this case, i.e. are the two probabilities the same? (Justify your answer)*

No, using normal tables, we obtain that $\mathrm{P}(X > 2) = \mathrm{P}((X - \mu)/\sigma > (2 - \mu)/\sigma)$ [1 mark], and $\mathrm{P}(Z > \sqrt{2}) = 1 - 0.9207 < 1/2$ [2 marks]

## Question 3 (17 marks in all)

The number of radioactive counts in 100 one minute intervals for a particular machine were

| No. of counts $(x)$: | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| Observed frequency: | 15 | 28 | 23 | 16 | 12 | 6 |

Suppose that we can model the number of counts in one minute by a Poisson random variable $X$, where

$$P(X = i) = e^{-\lambda}\lambda^i/i!, \quad i = 0, 1, 2, \ldots.$$

(a) (4 marks) *Prove that the probability generating function of the random variable $X$ is $\pi(s) = e^{\lambda(s-1)}$.* For $X \sim \mathcal{P}(\lambda)$,

$$
\begin{aligned}
\pi(s) &= \sum_{i=0}^{\infty} s^i \, \mathrm{P}(X = i) \\
&= \sum_{i=0}^{\infty} e^{-\lambda} \frac{\lambda^i}{i!} s^i \\
&= e^{-\lambda} \sum_{i=0}^{\infty} \frac{e^{\lambda s}}{e^{\lambda s}} \frac{(\lambda s)^i}{i!} \\
&= e^{-\lambda} e^{\lambda s} \\
&= e^{\lambda(s-1)}.
\end{aligned}
$$

(b) (2 marks) *Use the probability generating function in part (i) to prove that* $\mathrm{E}(X) = \lambda$.

Because, $\pi'(s) = \lambda e^{\lambda(s-1)}$ it follows that $\pi'(1) = \mathrm{E}\,X = \lambda$.

(c) (2 marks) *It is known that Var$(X) = \mathrm{E}(X) = \lambda$, determine* $\mathrm{E}(X^2)$.

From the definition of the $\mathrm{Var}(X) = \mathrm{E}(X^2) - (\mathrm{E}(X))^2$.

It follows that $\mathrm{E}(X^2) = \lambda + \lambda^2$.

(d) (2 marks) *Calculate the expected number of 0's in a sample of size 100 from a Poisson random variable with mean 2.*

$$\mathrm{P}(X = 0) = e^{-2} \quad \Rightarrow 100 \times e^{-2} = 13.53353.$$

(e) (7 marks) *You are given that the above sample average is $\bar{x} = 2$. Complete the table of expected frequencies below and test the goodness of fit of the Poisson distribution as a model for the number of radioactive counts.*

| No. of counts ($x$): | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| Expected frequency: | ? | 27.07 | 27.07 | 18.04 | 9.02 | ? |

The expected frequencies have to add to the sample size, i.e $E_5 = 100 - 9.02 - \ldots - 13.53353 = 5.26647$. [1 mark]

Thus the test statistic becomes,

$$X^2 = \sum O_i^2/E_i - n = \frac{15^2}{13.53353} + \ldots + \frac{6^2}{5.26647} - 100 = 2.12. \text{ [3 marks]}$$

Since here $g = 6$ and $k = 1$ the degrees of freedom is $\nu = g - k - 1 = 6 - 1 - 1 = 4$, we have that the $P$-value is larger than 0.25 since from the tables $P(\chi_4^2 > 5.385) = 0.25$. [2 marks]

Thus there is not sufficient statistical evidence to say that a Poisson distribution is not a good model. [1 mark]

## Question 4 (17 marks in all)

When cancerous tumours are removed from the colon it is not always possible to remove all cancerous cells without removing too much of the patient's vital organs. Consider the following data:

| | | Was the cancer controlled? | |
|---|---|---|---|
| | | Yes | No |
| Was cancer present | Yes | 8 | 182 |
| at the edge of surgery? | No | 11 | 58 |

(a) (8 marks) *Is there any evidence that cancer at the edge of surgery affects the chance of the cancer being controlled?*

This a $2 \times 2$ contingency table, i.e. with entries $x_{ij}$ the expected frequencies under independence are

$$E_{ij} = n \times \frac{x_{i\bullet}}{n} \times \frac{x_{\bullet j}}{n} = \frac{x_{i\bullet} \times x_{\bullet j}}{n},$$

where $x_{i\bullet}$ denotes the sum of the $i$th row and $x_{\bullet j}$ denotes the sum of the $j$th column. Here $n = 259$ and the row and column totals are:

|  |  | Was the cancer controlled? | | |
|---|---|---|---|---|
|  |  | Yes | No | Total |
| Was cancer present | Yes | 8 | 182 | 190 |
| at the edge of surgery? | No | 11 | 58 | 69 |
|  | Total | 19 | 240 | 259 |

Thus $E_{11} = \frac{190 \times 19}{259} = 13.93822$, $E_{12} = \frac{190 \times 240}{259} = 176.0618$, $E_{21} = \frac{69 \times 19}{259} = 5.061776$ and $E_{22} = \frac{69 \times 240}{259} = 63.93822$. [3 marks]

|                        |     | Was the cancer controlled? |     |
|------------------------|-----|:---:|:---:|
|                        |     | Yes | No  |
| Was cancer present     | Yes | 8   | 182 |
| at the edge of surgery? | No  | 11  | 58  |

From previous page $E_{11} = 13.93822$, $E_{12} = 176.0618$, $E_{21} = 5.061776$ and $E_{22} = 63.93822$.

Thus the test statistic becomes,

$$X^2 = \sum O_{ij}^2/E_{ij} - n = \frac{8^2}{13.94} + \ldots + \frac{58^2}{63.94} - 259 = 10.24814. \text{ [2 marks]}$$

The degrees of freedom here uses $g = 4$ and $k = 2$ so $g - k - 1 = 1$ or $(r-1)(c-1) = 1$ [1 mark].

The $P$-value is less than 0.01 since from the tables $\mathrm{P}(\chi_1^2 > 6.635) = 0.01$. [1 mark] Thus there is statistical evidence that that cancer at the edge of surgery affects the chance of the cancer being controlled. [1 mark]

(b) (5 marks) Consider the $69$ patients who had 'no cancer present at the edge of surgery' only:

*Provide a conservative $95\%$ confidence interval for the proportion of having controlled cancer.*

The formula for a conservative CI is $\widehat{p} \pm z_{\alpha/2}\sqrt{1/(4n)}$ [1 mark].

Here
$$\frac{11}{69} \pm 1.96/\sqrt{4 \times 69} = (\,0.0414\,,\,0.2774\,)\ \text{[2 marks]}$$

*Determine how much smaller the length of the conservative confidence interval in (i) is when the sample size was 100 instead of 69 patients.*

This is only determined by the $\sqrt{n}$ term, thus increasing the sample size from 69 to 100, makes the CI smaller by a factor of $\sqrt{100/69} = 1.2$. [2 marks]

(c) (4 marks) *Let $A$ and $B$ be two independent events. Show that* $\mathrm{P}(A \cap B^c) = \mathrm{P}(A) \times \mathrm{P}(B^c)$, *where $B^c$ denotes the complement of B.*

Note that $A = (A \cap B^c) \cup (A \cap B)$ and because those two events are mutually exclusive it follows from Axiom 3 that $\mathrm{P}(A) = \mathrm{P}(A \cap B^c) + \mathrm{P}(A \cap B)$ [1 mark].

Thus, $\mathrm{P}(A \cap B^c) = \mathrm{P}(A) - \mathrm{P}(A \cap B)$ [1 mark]

and using that $A$ and $B$ are independent we obtain

$$
\begin{aligned}
\mathrm{P}(A \cap B^c) &= \mathrm{P}(A) - \mathrm{P}(A)\,\mathrm{P}(B) \ [1 \text{ mark}] \\
&= \mathrm{P}(A)(1 - \mathrm{P}(B)) \\
&= \mathrm{P}(A)\,\mathrm{P}(B^c) \ [1 \text{ mark}]
\end{aligned}
$$

# Extended Answer Section of 2011 Exam

## Question 1 (14 marks in all)

The following R-output gives daily temperature, $x_i$ in degrees Fahrenheit, and Ozone level, $y_i$ in parts per billion in New York over 16 successive days.

```
> x = c(61,61,67,81,79,76,82,90,87,82,77,72,65,73,76,84)
> y = c(4,32,23,45,115,37,29,71,39,23,21,37,20,12,13,135)
```

Additionally you might find the following R output of use:

```
> sum(x)
[1] 1213
> sum(y)
[1] 656
> sum(x^2)
[1] 93125
> sum(y^2)
[1] 46868
> sum(x*y)
[1] 52111
> sort(round(lm(y~x)$resid,0))
  10   15    7    9   14   11    1    4    6    3    8   13   12    2    5   16
 -31  -28  -25  -25  -23  -22   -7   -7   -4    0    1    1    4   21   67   77
```

From these we can calculate the summary statistics:

$$S_{xx} = \left(\sum x_i^2\right) - \left(\sum x_i\right)^2/n = 93125 - 1213^2/16 = 1164.4375$$

$$S_{yy} = \left(\sum y_i^2\right) - \left(\sum y_i\right)^2/n = 46868 - 656^2/16 = 19.972$$

$$S_{xy} = \left(\sum x_i y_i\right) - \left(\sum x_i\right)\left(\sum y_i\right)^2/n = 52111 - 1213 \times 656/16 = 2378$$

# Part a (8 Marks)

## Subpart i

*Calculate the correlation coefficient. How would the correlation coefficient change if the daily temperatures were measured in degrees Celsius instead of degrees Fahrenheit?*

[1 mark for $r$, 1 mark for comment] The correlation coefficient is:

$$r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \frac{2378}{\sqrt{1164.4375 \times 19.972}} = 0.4931091$$

If the measurements where converted to Celsius from Fahrenheit then the correlation coefficient $r$ would remain unchanged.

# Subpart ii

*Calculate the least squares regression fit and the proportion variability of y's is explained by the regression line.*

[1 mark for $b$, 1 mark for $a$, 1 mark for $r^2$]

$$b = \frac{S_{xy}}{S_{xx}} = \frac{2378}{1164.4375} = 2.042187752$$

$$
\begin{aligned}
a &= \overline{y} = b\overline{x} \\
&= 656/16 - 2.042187752 \times 1213/16 \\
&= -113.8233589.
\end{aligned}
$$

The proportion ov variability explained by the regression fit is given by $r^2 = 0.4931091^2 = 0.243156584$.

# Subpart iii

*Use the* `R` *output above to calculate a 5 number summary. Use the 5 number summary to comment on the distribution of the residuals.*

[1/2 for min, 1/2 for max, 1/2 for median, 1/2 for $Q_1$, 1/2 for $Q_2$, 1/2 for $Q_3$ and 1/2 for comment] From the `R` output we have

- ☐ min $= -31$
- ☐ $Q_1 = \frac{-25-24}{2} = -24$
- ☐ $Q_2 = \frac{-7-4}{2} = -5.5$
- ☐ $Q_3 = \frac{1+4}{2} = 2.5$
- ☐ max $= 77$

From the 5 number summary we can see that the data is left skewed.

# Part b (6 Marks) − Subpart i

*Consider the paired data* $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$.

*For a least squares regression fit show that the residuals* $\widehat{e}_i = y_i - a - bx_i$ *satisfy*

$$\sum_{i=1}^{n} \widehat{e}_i = 0 \qquad \text{and} \qquad \sum_{i=1}^{n} \widehat{e}_i x_i = 0.$$

[1 mark for $\sum e_i = 0$, 3 marks for $\sum x_i \widehat{e}_i = 0$, 2 marks for $r = 0$] For the first part

$$\begin{aligned}
\sum e_i &= \sum (y_i - a - bx_i) \\
&= \sum n(\overline{y} - a - b\overline{x}) \\
&= 0
\end{aligned}$$

since $a = \overline{y} - b\overline{x}$.

Then

$$\sum x_i \widehat{e}_i = \sum (y_i - abx_i)x_i$$

$$= \sum (y_i - \overline{y} - b(x_i - \overline{x}))x_i$$

$$= \sum (y_i - \overline{y} - b(x_i - \overline{x}))(x_i - \overline{x}) + \overline{x}\sum (y_i - \overline{y} - b(x_i - \overline{x}))$$

$$= S_{xy} - bS_{xx} + 0 \quad (\text{since } \sum (y_i - \overline{y}) = \sum (x_i - \overline{x}) = 0)$$

$$= S_{xy} - \frac{S_{xy}}{S_{xx}}S_{xx}$$

$$= 0$$

# Subpart ii

*Hence, show that $\{x_i\}$ and $\{\widehat{e}_i\}$ are uncorrelated.*

To calculate the (sample) correlation of $x_i$ and $\widehat{e}_i$ we have

$$r = \frac{S_{xe}}{\sqrt{S_{ee}S_{xx}}}$$

where

$$S_{xe} = \left(\sum x_i\widehat{e}_i\right) - \left(\sum x_i\right)\left(\sum \widehat{e}_i\right)/n = 0 - \left(\sum x_i\right) \times 0/n = 0$$

since $\sum e_i = 0$ and $\sum x_i\widehat{e}_i = 0$. Hence $r = 0$.

# Question 2 (18 marks in all)

## Part a (9 Marks)

*The clinically accepted value for mean blood pressure in healthy males aged 18 to 22 years is 120 mm Hg. It is widely claimed that examination stress causes blood pressure to rise above 120 mm Hg. To test this theory, 10 healthy male students have their blood pressure taken just prior to a Statistics quiz. The sample mean and sample standard deviation of these measurements are 135.1 and 17.42 respectively. Assume that the measurements for each student are independent and normally distributed.*

Note: We have from the description $n = 10$, $\overline{x} = 135.1$ and $s = 17.42$.

# Subpart i

*What are appropriate null and alternative hypotheses to test this claim?*

[1 mark for $H_0$, 1 mark for justifying $H_1$] Since "it is widely believed that exams increase blood pressure" the appropriate null and alternative hypotheses are:

$$H_0\colon \mu = 120 \qquad \text{vs} \qquad H_1\colon \mu > 120.$$

# Subpart ii

*State an appropriate test statistic to test this hypothesis and the null distribution of this test statistic.*

[1 mark for $T$, 1 mark for justifying $t_9$] The appropriate test statistic is

$$T = \frac{\overline{x} - \mu_0}{s/\sqrt{n}}$$

where $\mu_0 = 120$. Assuming independence and normality under the null hypothesis $T \sim t_{n-1} = t_9$.

# Subpart iii

*Calculate the test statistic chosen in (ii), the corresponding $P$-value and form an appropriate conclusion.*

[2 mark for $T$, 1 mark for $P = P(t_9 > 0.274112472)$, 1 mark for $P > 0.25$, 1 mark for retaining null, 1 mark for relating to relating to question] The value of the test statistic is

$$T = \frac{135.1 - 120}{17.42/\sqrt{10}} = 0.274112472$$

The P-value is given by

$$P = P(t_9 > 0.274112472) > 0.25 \qquad \text{(via tables)}.$$

Hence we retain the null hypothesis. There is insufficient evidence to suggest that exam stress causes blood pressure to rise above 120mm Hg.

# Part b (9 Marks) – Subpart i

*Consider the geometric distribution $P(X = x) = pq^x$, $x = 0, 1, 2, \ldots$, $0 \leq p \leq 1$ and $q = 1 - p$. Show that the probability generating function is given by $\pi(s) = p/(1 - qs)$ for $|s| < 1/q$.*

[1/2 mark for defining $\pi(s)$, 1/2 for $p\sum_{i=0}^{\infty}(sq)^i$, 1/2 mark for $\frac{p}{1-qs}$, 1/2 mark for $|s| < 1/q$] We have

$$P(X = x) = pq^x, \qquad x = 0, 1, 2, \ldots$$

The probability generating function is

$$\pi(s) = \sum_{i=0}^{\infty} s^i P(X - i) = \sum_{i=0}^{\infty} s^i pq^i = p\sum_{i=0}^{\infty}(sq)^i = \frac{p}{1 - qs}$$

provided $|qs| < 1$ or equivalently $|s| < 1/q$ (since $q > 0$).

# Subpart ii

*Use Part (i) to show that $E(X) = q/p$ and Var$(X) = q/(p^2)$.*

[1 mark for $\pi'(s)$, 1 mark for $\pi''(s)$, 1 mark for $E(X)$, 1 mark for Var$(X)$] Note that

$$\pi'(s) = \frac{pq}{(1 - qs)^2} \qquad \text{and} \qquad \pi''(s) = \frac{2pq^2}{(1 - qs)^3}.$$

Now

$$E(X) = \pi'(1) = \frac{pq}{(1 - q)^2} = \frac{pq}{p^2} = \frac{q}{p}$$

and

$$\text{Var}(X) = \pi''(1) + \pi'(1) - \pi'(1)^2 = \frac{2pq^2}{(1 - q)^3} + \frac{p}{q} - \frac{q^2}{p^2} = \frac{q(p + q)}{p^2} = \frac{q}{p^2}.$$

# Subpart iii

*Use Part (ii) and Chebyshev's inequality to bound $P(|X - q/p| > 1)$.*

[1/2 mark for Chebyshev's inequality, 1/2 mark for identifying $\sigma$, 1 for identifying $c$, 1 mark for boudn] Chebyshev's inequality is

$$P(|X - \mu| > c\sigma) < 1/c^2.$$

We have $\mu = q/p$ and $\sigma = \sqrt{q}/p$.

$$P(|X - p/q| > 1) = P(|X - q/p| > p/\sqrt{q} \times \sqrt{q}/p) < 1/(p/\sqrt{q})^2 = q/p^2$$

where $c = p/\sqrt{q}$.

# Question 3 (15 marks in all)

## Part a (10 marks)

*It has been claimed that at least 60% of all purchasers of a certain computer program will call the manufacturer's hotline within one month of purchase. A random sample of 12 purchasers of this software is drawn and 3 of those in the sample had contacted the hotline within one month of purchase. Does this provide evidence that the claim of a 60% contact rate is an overestimate? Let $p$ be the true proportion of all purchasers who contact the hotline.*

From the question description we have $p_0 = 0.6$, $n = 12$ and $X = 3$.

# Subpart i

*Calculate an approximate 95% confidence interval for $p$.*

[1 mark for formula, 1 mark for $\widehat{p}$, 1 mark for calculation] The 95% confidence interval for $p$ is given by

$$\widehat{p} \pm 1.96 \sqrt{\frac{\widehat{p}(1-\widehat{p})}{n}} = \frac{1}{4} \pm 1.96 \sqrt{\frac{3/16}{12}} = 0.25 \pm 0.245 = (0.005, 0.495).$$

where $\widehat{p} = 3/12 = 1/4$.

## Subpart ii

*Form an appropriate hypotheses and perform a statistical test for the above situation stating any assumptions you may require.*

[1 mark for assuming $X \sim B(n, p)$, 1 mark for hypotheses, 1 mark for null distribution, 1 mark for P-value, 1 mark for conclusion] Assume constant proability over each trial, each trial is independent so that $X \sim B(n, p)$. Appropriate null and alternative hypotheses are:

$$H_0\colon p = 0.6 \qquad \text{and} \qquad H_1\colon p < 0.6.$$

The test statistic is $X = 3$ with small values of $X$ supporting $H_1$. Under $H_0$ we have $X \sim B(12, 0.6)$. Then

$$\begin{aligned} P(X \leq 3) &= P(X = 3) + P(X = 2) + P(X = 1) + P(X = 0) \\ &\approx 0.01246 + 0.0025 + 0.0003 + 10^{-5} \\ &\approx 0.0153 \end{aligned}$$

Also, using the normal approximation $X \approx Y \sim N(np, np(1-p)) = N(7.2, 2.88)$ and

$$
\begin{aligned}
P(X \leq 3) &\approx P(Y \leq 3.5) \\
&= P(Z \leq (3.5 - 7.2)/\sqrt{(2.88)}) \\
&= P(Z \leq -2.180246) \\
&= 1 - P(Z \leq 2.18) \\
&= 1 - 0.9854 \\
&= 0.0146
\end{aligned}
$$

(from tables) where $Z \sim N(0, 1)$. Either way the P-value is less than 5%. Hence we reject $H_0$ and conclude that there is sufficient to suggest that 60% is too high.

## Subpart iii

*Do the results from part (i) and part (ii) agree? Justify your answer.*

[1 mark for stating CI does not contain $0.6$, 1 mark for comment] Since that 95% CI does not contain the value $p_0 = 0.6$ and we reject $H_0\colon p = 0.6$ the conclusions are consistent.

## Part b (5 marks)

*Suppose that the probability that a randomly chosen individual having a particular disease is 0.02.*

From the problem description we have $n = 100$ and $p = 0.02$ with $X \sim B(100, 0.02)$.

## Subpart i

*Write an expression for the exact probability that at most 2 cases in a random sample of 100 has the disease.*

[1/2 mark for $n$, 1/2 mark for $p$, 1 mark for $P(X \leq 2)$] The expresssion for $P(X \leq 2)$ is

$$P(X \leq 2) = \binom{100}{2}(0.02)^2(0.98)^{98} + 100(0.02)(0.98)^{99} + (0.98)^{100}$$

## Subpart ii

*Use the Poisson approximation to the Binomial distribution to approximate this probability. Why is this approximation appropriate?*

[1 mark for $Y \sim P(2)$, 1 mark for $P(Y \leq 2) = 5e^{-2}$, 1 mark for comment] The Poisson approximation is $X \approx Y \sim P(np) = P(2)$. Then

$$P(Y \leq 2) = \frac{2^2 e^{-2}}{2} + \frac{2e^{-2}}{1} + \frac{2^0 e^{-2}}{1} = 5e^{-2} = 0.6766764.$$

This approximation is appropriate since $n = 100$ is large and $p = 0.02$ is small.

## Question 4 (18 marks in all) – Part a (5 marks)

*A new casino game involves rolling 3 dice. The winnings are directly proportional to the total number of sixes rolled. Suppose a gambler plays the game 101 times, with the following observed counts:*

| Sixes Rolled | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| Number of Rolls | 48 | 35 | 15 | 3 |

*You may use the R output:*

```
> dbinom(0:3,3,1/6)
[1] 0.57870370 0.34722222 0.06944444 0.00462963
```

*The casino becomes suspicious of the gambler and wishes to determine whether the dice are fair. Use a statistical test to determine where whether the dice are fair or not.*

[1/2 mark for $g$, 1/2 mark for $E_i$, 1/2 mark for combining, 1 mark for $X^2$, 1/2 mark for $X^2 \sim \chi_2^2$, 1/2 mark for P-value, 1/2 mark for conclusion] There are $g = 4$ groups and

| $O_i$ | 48 | 35 | 15 | 3 |
|-------|------|--------|-------|-------|
| $E_i$ | 58.449 | 35.069 | 7.014 | 0.468 |

Since the last column contains both expected and observed frequencies less than 5 we combine the last two columns to obtain the table:

| $O_i$ | 48 | 35 | 18 |
|-------|------|--------|-------|
| $E_i$ | 58.449 | 35.069 | 7.481 |

The observed value of the Peason's chi-square statistic is

$$X^2 = \sum_{i=1}^{g} \frac{O_i^2}{E_i} - n = 38.81 + 35.28 + 43.74 - 101 = 16.66 \quad \text{(to 2 d.p.)}$$

Here $g = 3$. Under the null hypothesis $X^2 \sim \chi_{g-1}^2 = \chi_2^2$ and the P-value is given by $P(\chi_2^2 > 16.66) < 0.01$ from tables. This we reject the claim that the dice are fair.

## Part b (9 marks)

*Two pathology labs, lab A and lab B, are compared to see which of the labs report their results for a specific test faster. Samples from 20 patients are collected and then 10 of these samples are sent to lab A and 10 of these samples are sent to lab B. Summary values for the times, in days, for each lab to report its results are summarised in the table below:*

|       | Size | Mean  | Median | Variance |
|-------|------|-------|--------|----------|
| Lab A | 10   | 20.23 | 19.45  | 2.74     |
| Lab B | 10   | 18.68 | 17.98  | 1.64     |

*You may assume that the measurements for lab A and lab B are normally distributed. Suppose we wish to test whether there is a difference in the times each lab reports its results.*

# Subpart i

*State an appropriate null and alternative hypothesis, defining any parameters used.*

[1/2 mark for $H_0$, 1/2 mark for $H_1$, 1 mark for definitions]

$$H_0 \colon \mu_x = \mu_y \qquad \text{vs} \qquad H_1 \colon \mu_x \neq \mu_y$$

where $\mu_x = \{\text{mean time form lab A}\}$ and $\mu_y = \{\text{mean time form lab B}\}$.

## Subpart ii

*State an appropriate test statistic to test this hypothesis and the null distribution of this test statistic, stating any additional assumptions required.*

[1 mark for justifying 2 sample t-test, $1/2$ mark for test statistic, $1/2$ mark for null distribution, 1 mark for equal variance assumption] The two sets of samples are independent so that a 2 sample t-test is appropriate. The appropriate test statistic is:

$$T = \frac{\overline{X} - \overline{Y}}{S_p\sqrt{1/n_x + 1/n_y}}$$

Assuming equal variance the null distribution of $T$ is $t_{n_x+n_y-2} = t_{18}$.

## Subpart ii

*Calculate the test statistic chosen in (ii), the corresponding $P$-value and form an appropriate conclusion.*

[1 mark for $s_p^2$, 1 mark for $T$, 1 mark for P-value, 1 mark for conclusion] The pooled variance is

$$s_p^2 = \frac{(n_x - 1)s_x^2 + (n_y - 1)s_y^2}{n_x + n_y - 2} = \frac{1}{2}(2.74 + 1.64) = 2.19$$

The observed value of the test statistic is

$$T = \frac{20.23 - 18.68}{\sqrt{2.19}\sqrt{2/10}} = 2.342.$$

The P-value is given by

$$P = P(|t_{18}| > 2.2342) = 2P(t_{18} > 2.2342).$$

From tables $P(t_{10} > 2.228) = 0.025$ and $P(t_{20} > 2.528) = 0.01$. Hence $0.02 < P < 0.05$. We reject $H_0$ and conclude that the two labs have different mean times.

## Part b (4 marks)

Let $\Omega$ be a sample space, $A \subset \Omega$, $B \subset \Omega$ and $P(\cdot)$ be a probability function satisfying the axioms of probability

☐ For any event $A \subset \Omega$, $P(A) \geq 0$,

☐ $P(\Omega) = 1$

☐ If $A$ and $B$ are mutually exclusive events $(A \cap B = \emptyset)$, then

$$\mathrm{P}(A \cup B) = \mathrm{P}(A) + \mathrm{P}(B).$$

$$
\begin{aligned}
1 &= P(\Omega) && \text{by axiom 2 (1/2 mark)}\\
P(\Omega) &= P(A \cup A^c) && \text{by set properties (1 mark)}\\
P(A \cup A^c) &= P(A) + P(A^c) && \text{by axiom 3 (1 mark)}\\
P(A) &= 1 - P(A^c) && \text{(1/2 mark)}\\
P(A) &= 1 - P(A^c) \leq 1 && \text{since } P(A^c) \geq 0 \text{ by axiom 1 (1 mark)}
\end{aligned}
$$