

93201A



S

12.30

SUPERVISOR'S USE ONLY

SCHOLARSHIP EXEMPLAR



NEW ZEALAND QUALIFICATIONS AUTHORITY
MANA TOHU MĀTAURANGA O AOTEAROA

Scholarship 2014 Statistics

9.30 am Wednesday 12 November 2014

Time allowed: Three hours

Total marks: 40

ANSWER BOOKLET

Check that the National Student Number (NSN) on your admission slip is the same as the number at the top of this page.

Write all your answers in this booklet.

Show ALL working. Start your answer to each question on a new page. Clearly number each question.

Check that this booklet has pages 2–24 in the correct order and that none of these pages is blank.

YOU MUST HAND THIS BOOKLET TO THE SUPERVISOR AT THE END OF THE EXAMINATION.

1) a) The Diabetic rate vs Obesity rate graph shows a weak positive correlation between the two variables. As the percentage of the adult population who are obese increases, the diabetes rate also increases. For every 1% rise in the obesity rate, the Diabetic rate tends to rise by 0.0344%. The ~~graph~~ regression line intercepts the y axis at 2.3718%. I expect that this is due to the proportion of the population who have diabetes as a result of their genetic makeup, ~~which will be constant across~~ This will be approximately constant across the ~~countries included~~ countries included, and is ~~not affected by~~ not included in the relationship between a correlation relationship with the Obesity rate of a country. The diabetic rate vs GDP per capita establishes a moderate positive relationship, where as the GDP increases, the Diabetic rate also tends to increase. For every \$10 000 increase in GDP per capita, the Diabetic rate tends to increase by 0.6445%. It is interesting to note that there is far more variation towards the lower end ~~and around~~ in the Diabetic rate for countries with low a GDP. This may just be a feature of our sample however, as the distribution of GDPs in our sample is significantly in the favour of low GDP values, which may generate a perception in the graph that there is more variation around ~~the~~ low GDPs. The Life expectancy vs Diabetic rate shows the strongest ^{positive} relationship of all the three graphs, ~~that~~ It shows us that as Diabetic rate increases, the Life expectancy also increases. ~~the~~ A 1% increase in the Diabetic rate ~~will tend to~~ would tend to mean a ~~1.0671 year~~ increase in the Life

It may be possible to generate a more accurate model using multivariate analysis to find an relationship between both drug₁ rate and GDP, and hence they affect the numbers rate.

Using GDP as an indicator:

$$y = 0.6445(0.5) + 2.1002$$

$$\approx 2.4224596 \approx 2.4910 (2sf)$$

Using obesity rate as an indicator

$$y = 0.0344(71) + 2.3718$$

$$\frac{2.4424}{2} - 0.371P$$

$$\approx 4.81429, \approx 4.8\% (28f)$$

- I would be much more confident using the ~~Estimation~~ regression estimated using the GDP, as the ~~relationship between~~ the GDP and default rate is stronger than the relationship between ~~Unemployment rate and default rate given in Equation 2~~

Fig 1 & 2. The stronger relationship means I am ~~more~~ confident that the actual value will be closer to the estimate when I use the GDP to estimate than when I use the Obdys rate.

I am reasonably confident ~~that this is correct~~, as these two assumptions, as the two y intercepts on the two decades graphs are within an acceptable distance of each other (2378 in Fig 1 compared to 21002 in Fig 2). If I am correct in this assumption then there would be close to identical as either the density rate or GDP would have any effect on the position of the population with densities

- 4.1) Clearly, my ~~estimated~~ estimates of 2.49% and 1.87% were drastically off from the actual value of 2.2%. Therefore there must be significant confounding variables that were ~~been~~ a major exception from the trends ~~of the~~ identified in the graphs. Looking at the information presented, I expect the reason to be the poor state of New Zealand's healthcare. The poor condition of a country's healthcare system would mean that ~~there~~ both the ~~overseas~~ prevention and treatment of diabetes patients and ~~the~~ education of those who are at risk of developing diabetes due to their lifestyle habits would suffer. This could be the reason for the ~~diabetes~~ significantly higher diabetes rate than the model's predictions.
- 4.2) With 95% confidence, I can say that ~~the~~ in 2006/7 the Obesity rate ~~could have been~~ for New Zealanders could have been as low as 25.5% or as high as 27.5%, and that in 2012/13 the ~~the~~ obesity rate could have been as low as 27.4% or as high as 29.8%. As there is overlap between these two ranges, ~~into the~~ ~~obesity rate~~ it is possible that the obesity rate dropped from 27.5% in 2006/7 to 27.4% in 2012/13, ~~the~~ ~~claim~~ that the obesity rate is decreasing is unsupported.

2) (i) Random assignment was used when forming the groups to ~~prevent~~ minimise the effect of non-sampling errors such as self-selection bias. ^{For example} If the students were able to choose what group they were placed in, then those who were wanting to lose weight may have chosen to go into the group that did ~~nothing~~ ^{lose weight}. ~~By random assignment~~ ^{the results would have been affected by this, as self-selecting students} leaves a lot more variables uncontrolled, which may have ~~been~~ ^{not} randomly assigning students leaves a lot more variables uncontrolled (close to lose weight, for example), which may have a confounding effect on the results obtained.

(ii) From this output, I can conclude that there is a 8.4% probability that the difference visible at the end of the experiment was due to chance acting alone. The randomisation process ~~was~~ randomly assigns ~~the~~ ^{the} results of each participant into one participant ~~consequence~~ many times. The purpose of this is to emulate the variability within the population, and to find out how likely it would have been to often ~~conclude~~ ^{conclude} the same, or greater difference if we had directed our two groups differently. The 8.4% result generated ~~concerns~~ ^{is significant} large enough to make me question whether the difference visible is due to chance acting alone, ~~however~~ ^{or} due to the different treatment of the groups. I am not confident that the claim that the activities prevented ~~weight gain~~ ^{weight gain} can help children stop being obese is supported by this study.

41.) The two surveys found that within the sample of people who took the survey, there was a drop in the mean number of sugary drinks consumed of 0.6 (4.9-4.3). Interestingly, the second survey results exhibited less variation than the first. This is reflected in the drop in the standard deviation from 2.6 to 2.1, and the decrease in the ~~sample~~ range covered by the middle 50% of ~~the samples~~ from 4 in the first survey (7-3), to 2.5 in the second survey (5.5-3). Both graphs show the distributions to be approximately normal, although the ~~first~~ although the second survey has a normal distribution worse than the first, as it ~~is~~ is skewed slightly towards higher values.

42.) The purpose of the bootstrap distribution is to emulate what would happen if we took a very large number (usually 1000 using iMight) of samples from the population. It does this by randomly resampling ~~the sample groups~~ from the original two ~~participate from~~ the two survey response groups from the existing results obtained with replacement, if the ~~same~~ person's response ~~is~~ may be present in the resampled sample multiple times, whilst others may not be present at all. The ~~distribution~~ is then plotted on the ~~bootstrap distribution~~. This process is then repeated many times to produce the distribution shown. Finally the top and bottom 2.5% of the ~~points on the bootstrap distribution~~ are removed and

~~removed~~ and we note. Firstly, we take the range covered ^{between} the 2.5th percentile and the 97.5th percentile, and use 1308 for the range in which we can be 95% confident that the ~~difference~~ population difference in the ~~rate of sugary drinks consumed~~ as a ~~result~~ ~~of the~~ is the ~~number~~ number of sugary drinks consumed prior to, and following the education campaign lies between. Using the bootstrap distribution, I can be 95% confident that on average 112 secondary school students would drink between 0.10 and 1.24 sugary drinks less following the ~~ad~~ education and advertising campaign. Therefore, as 0 is not within this range, ~~we can~~ ~~conclude~~ ~~conclude~~ ~~in support~~ the claim that the campaign made 112 secondary school students drink less sugary drinks is supported. //

Q.2.
PSP
5

3) a) Figure 8 shows a clear upward trend in the obesity rate of both men and women in New Zealand. ~~From~~ The percentage of men who were obese grew from 24.5% in 2004 to 29.0% in 2013, at an average rate of 0.39% per year.

$$\text{av rate of growth} = \frac{\text{final} - \text{initial}}{\# \text{ time periods}}$$

$$= \frac{29 - 24.5}{9}$$

$$= 0.388889$$

$$= 0.39\% (2sf)$$

The male line shows a far greater consistency in the rate of increase of the obesity rate than the female, as evidenced by the much straighter line. The female line also shows a clear upward trend, although it is much more erratic, and ~~also~~ ~~relates~~ to the point where it indicates drops in the obesity rate in the 2006-07, 2008-09 and 2011-12 periods. The female obesity rate climbed from 26.4% to an 29.5% over the period from 2004 to 2013, at an average rate of 0.34%

$$\text{av rate of growth} = \frac{\text{final} - \text{initial}}{\# \text{ time periods}}$$

$$= \frac{29.5 - 26.4}{9}$$

$$= 0.3444$$

$$= 0.34\% (2sf)$$

The female obesity rate is consistently higher than the male obesity rate across this period //

2004-2013

Start each question on a new page.

Figure 9 shows an increasing trend in the Obesity rates of ^{with} 15-24 year olds from 12.6% to 22.3% at an average rate of ~~1.19%~~ 1.19% per year and 35-44 year olds from 27.6% to 32.8% at an average rate of 0.58% per year, ~~and the 55-64 year olds~~. It also shows a decreasing trend in the Obesity rate of 55-64 year olds from 37.1% in 2004 to 31.4% in 2013, at an average rate of -0.58% per year. The 15-24 and 55-64 year old age groups show a fairly consistent rate of change across this period, whilst the trend shown for the 35-44 age group does not follow an approximately linear path. From 2004-2008, the Obesity rate for 35-44 year olds rapidly ~~dropped~~ at an increasing rate till it reaches its maximum at 33.9% in 2008. In the next period the obesity rate dropped to 32.5% and remained relatively constant through till the end of the time period //

- 4) Using my graphics calculator I plotted the data points for the ~~massive~~ male Obesity rates across this period and used the calculator to generate a linear regression line that I ~~will~~ use to model the trend. The model equation generated was //

$$y = 0.3861x + 25.5$$

where y is the obesity rate and x is years since 2004. This line had a correlation coefficient of 0.9999 (3sf) so I can be ~~the~~ confident that it is a good tool for

predicting the obesity rate in 2015.

Using this model to predict the male obesity rate in 2015

$$y = 0.3861(10) + 29.5$$

$$= 29.4\% \text{ (3sf)}$$

I used a similar method to generate a regression line for the female obesity rate. ~~transposed the data~~
This regression line had the following equation:

$$y = 0.3061x + 26.54$$

This line had a correlation coefficient with the time series period of 0.924. Whilst this still means that the model generated will be reasonably accurate, the ~~slightly~~ less consistent nature of the female trend means that I will be less confident with this estimate than I was with the one for the males.

Using the equation to predict the ~~male~~ female obesity rate in 2015:

$$y = 0.3061(10) + 26.54$$

$$= 29.6 \text{ (3sf)}$$

- c) I notice that the female obesity rate tends to ~~decrease~~ drop every 3 years (i.e. in 2006, 2009, 2012). The cyclic nature of this must mean that there is some confounding variable ~~which~~ affecting female obesity rates which acts on a 3 year cycle.

ii) The difference in the two samples (assuming ideal experimental design) is due to sampling variation. When ever we take a survey, we are looking only at a small proportion of the total population. These two surveys would have taken different samples, from the same population, and whilst they would be from the same population, it is highly unlikely that they would have produced the same results. The 0.2% difference in the two conclusions is due to the two samples each containing different members of the same population. This difference would be well within the range for sampling error, and therefore it is not a concern. This is assuming the experiments were ideally designed however, as there may be some non sampling error in one of the experiment surveys that is partially to blame for the difference of 0.2%.

See Barber information on this //

Q.3
SNP
5

41 Chase not

65-74 0.386

$$\# \text{ people in } 65-74 \text{ age group who are Chase} = 6.14 \times 117000 = 196100$$

$$\# \text{ people in } 65-74 \text{ age group} = \frac{196100}{0.386} \times 100$$

$$= 404404$$

$$\approx 404,000 \text{ (3 sf)}$$

Chase	not	
156100		404,000

$$\text{not Chase} / 65-74 = \frac{156100}{404000}$$

a) $X \sim \text{Binomial}(1000, 0.386)$

I am using Binomial as: we can

• we can ~~take the 1000 people as a trial~~

• we can treat each of the 1000 people as a trial

• 1 person is either Chase or not

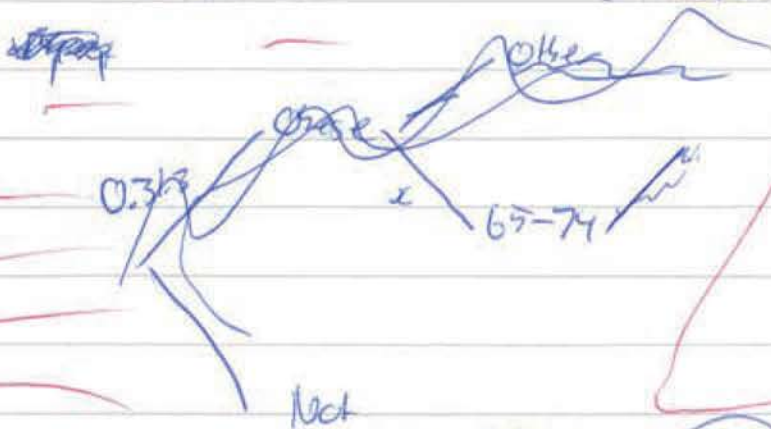
• The probability of one person in our randomly chosen sample being Chase is not affected by whether any other person is Chase

• There is a constant probability of someone being Chase of 0.386

$$P(X \leq 360) = 0.04836 \text{ (4 sf)}$$

Start each question on a new page.

h) # people in 65-74 age group = $0.14 \times 1115,000$
 & obese = 156100



$$P(\text{Obese} \cap 65-74) = 0.313 \times$$

$$P(\text{Obese} | 65-74)$$

$$P(65-74 | \text{Obese}) = \frac{P(65-74 \cap \text{Obese})}{P(\text{Obese})}$$

people in 65-74 age group
 $= \frac{156100}{386} \times 100$
 $= 404,000 (3sf)$

$\frac{0.313}{0.313}$

people over 15 = $\frac{1115000}{31.3} \times 100$
 $= 3,562,300$
 $= 3,562,000 (3sf)$

oo $P(65-74 \cap \text{Obese}) = \frac{156100}{3,562,000}$
 $= 0.04382$

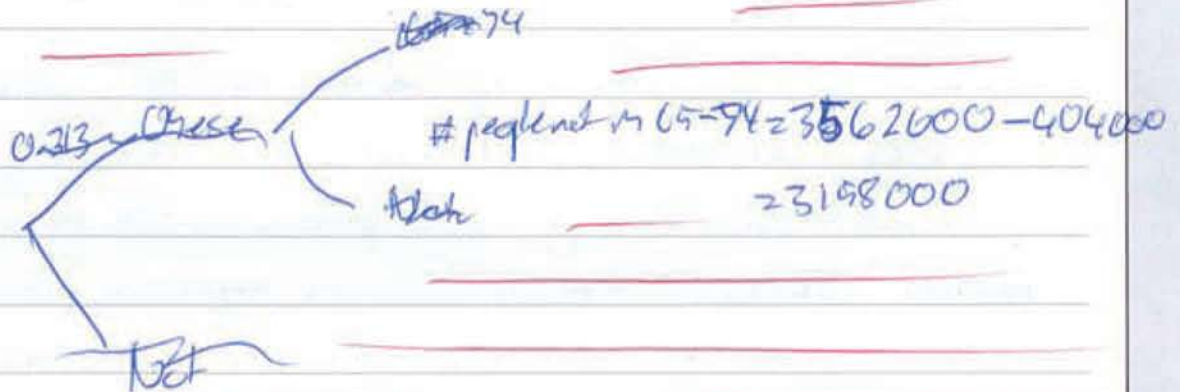
$P(65-74 | \text{Obese}) = \frac{P(65-74 \cap \text{Obese})}{P(\text{Obese})}$

$= \frac{0.04382}{0.313} = 0.14$

Start each question on a new page.

c) # obese people in 65-74 = 156100

obese not 65-74 = 1,150,000 - 156100
= 993900



$P(\text{obese given not in 65-74}) = \frac{834000}{3158000}$
= 0.3036 (4sf)

$\frac{P(\text{obese given 65-74})}{P(\text{obese given not 65-74})} = \frac{0.386}{0.3036}$

= 1.311

A person in the 65-74 age group is 1.3 times as likely to be obese as someone who is not within this age range.

Q.4.
NPO

7

QUESTION
NUMBER

Start each question on a new page.

ASSESSOR'S
USE ONLY

5) c) i) In Brazil, the prevalence of overweight children (6-9) increased from 4.9 to 17.4 in our samples, an increase of 12.5, although this increase could be as much as 14.03 or as low as 10.49 in the total population of Brazilians aged between 6-9. This ~~was an average rate~~ of sample growth was 0.54% per year, although this could be as low as 0.47% or as high as 0.61% in the population. This ~~is~~ trend shown in Brazil is notably different to the trend in overweight ~~the~~ prevalence in the same age group in Russia, where ~~the~~ which ~~appears~~ decreased by 16.2 ~~the~~ percentage points across the period from 1992-1999, although this change could be as low as 13.84 or as high as 18.56. Whilst we can make a comparison across these two countries, there is a large potential for non-sampling error. This is because the time span the information ^{was} taken across is far longer in ~~the~~ Brazil than it is in Russia, whilst the end dates are similar, the starting date of 1974 for Brazil and 1992 for Russia ~~are~~ are far too far apart for any meaningful insight to be gained. The Brazil values may have experienced significant change due to ~~the~~ some ^{global} event that would not have affected the data from Russia as they were yet to start recording their data.

In Russia, the prevalence of obesity amongst adolescents ~~was~~ ~~has~~ dropped from ~~the~~ 11.5 ~~to~~ to 8.5 across this ~~period~~ period, although this ~~drop~~ ~~is~~ ~~not~~ whilst

The prevalence of obesity amongst adolescents in China rose from ~~19.5%~~ ^{4.5%} ~~19.5%~~ ^{19.5%}. The Russian data shows a clear drop, as there is no overlap in the confidence intervals. However, the Chinese data ~~shows only a small increase~~ ^{shows only a small increase}, which could be as ~~low as~~ ^{low as} 0.68 or as high as 7.22 percentage points. The change in Russia ~~is a result of a number of factors~~ ^{is a result of a number of factors} suggests confounding the change of one or more factors which are strongly influencing the prevalence of obesity amongst adolescents. The author suggests this to be the period of tremendous economic stress the country went through.

It is interesting to note that in America, it is the lowest income households that ~~show~~ ^{have} the highest prevalence of obesity. This ~~is different to the~~ ^{is different to the} contrasts with Brazil, whose highest income families ~~have~~ ^{have} a higher prevalence of obesity, ~~as seen in~~ ^{as seen in} the 1997, (although their 1974 results show the lower income households to have a higher prevalence of obesity), and China and Russia, who have relatively even distributions of obesity across the three income levels measured.

- ii) The prevalence of underweight children in Brazil dropped significantly from 13.3 to 6.9% across the period from 1974 to 1997, at an average rate of 0.279% per year. ~~Whereas~~ ^{Whereas} whilst this

QUESTION
NUMBER

Start each question on a new page.

ASSESSOR'S
USE ONLY

re the largest drop //

h) The ~~worldly~~ ~~comparative~~ ~~significant~~ variables in the time periods covered, across the ~~three~~ ~~for~~ countries mean that it is difficult to make any comparisons between them. My suggestion would be to adapt the title to include ~~and~~ ^{country} rate of change per year column, as this would be very useful to ~~judge~~ ^{judge} compare trends across the ~~countries~~ countries. This would ~~minimise~~ ^{minimise} the influence of the confounding variable of the time ~~that~~ between the two surveys in each country, by allowing researchers and analysts to clearly ~~see~~ ^{compare} the rate at which obesity prevalence is changing across these three countries.

My suggestion for ~~the~~ Figure 1 would be to include ~~the~~ information on the prevalence of underweightness. Whilst the figure is a good tool to see how income affects overweightness across these three countries, it provides no information on the underweight, which is needed to show changes.

2P

c) The question is asking me to look at the changes over and under weight prevalence. I would expect that an analysis of the changes in these variables would include an analysis of the rate at which they are changing, especially, when the absolute values are difficult to compare across the countries as they are measured across sequentially different time periods //

N

The reasoning behind this suggestion is to help the viewer identify changes in the prevalence of underweight, which is not required to support the argument. The table currently provides no information that can be used to argue any thing about changes in underweight. /

N

Q4.
2p
2

Scholarship

1. Three points were correct from the first two graphs so an S was scored.
2. The spurious correlation was not picked up due to the fact that the obesity rate was correlated with both life expectancy and the diabetic rate.
3. Three points were picked up namely the best prediction value, an appropriate justification as to why it was best along with why Nauru could be regarded as an outlier with its very high obesity rate.
4. There was no mention of the difference in the point estimates which would have scored an S.
5. There was no mention of the actual difference along with the tail probability so answer was incomplete so a P was scored.
6. All three distinct comparisons were made; the drop in mean, the reduction in variation and a comparison of distributions.
7. There was no mention that we were looking at a range pertaining to the mean number of sugary drinks so a P was scored.
8. The prediction for the male obesity rate was incorrect and no account was taken of the fluctuations in the women's obesity rate in the other forecast calculation.
9. This question part was asking for the percentage change in the overall obesity rates and wasn't a repetition of part (a)'s requirements.
10. Only two points earning credit were made; a correct trend comparison between Brazil and Russia for overweight and a comment on how the data presentation in the table could be improved.