# STATS SAMPLE PAPER.

1. a Categorical data has no order or no quantity involved, that rules out alot of the statistical tests we run any at gender/days of week

Quantile data had an underlying scale and a quantity involved - must easier to work with this data
eg height/salary

b - Boxplot - shows three main features of each variable; its center represented by the horizontal line in the box, its spread - the vertical line and its outliers represented by stars
- Y axy measures the frequency
- x axis represent variable yes and no

c - The mean is commonly known as the average.
- Calculated by adding up all quantities and dividing by the number of quantities
- It is a measure of centre
- Can be affected by outliers

- Also known as the midpoint
- Median is the middle score of a list of scores
- less sensitive to extreme values (outliers)
- It is the point at which half the values are above and half below

D Range
- Difference between lowest and highest value in dataset
- Simple to compute
- Can also split up into interquartile range to measure the dispersion based upon two values from the dataset

Standard dev

- More powerful measure of variability because it takes into account every value in the data
- Measure of amount by which every value within a dataset varied from the mean
- Square root of variance

1e    $\bar{x} = 141$        $\sigma = 2.5$

$$\frac{141-137}{2.5} \qquad \frac{141-144}{2.5}$$

$$\frac{8}{5} \ (1.6) \qquad -\frac{6}{5} \ (-1.2)$$

$$-1.2 \leq z \leq 1.6$$

$z \leq 1.6$                    $-1.2 \leq z$
$= 0.945$                      $-z \leq 1.2$
$-0.12$                        $1 - (z \leq 1.2)$
$= 82.5\%$                     $1 - 0.88 = 0.12$

3

# STATS Sample Paper.
## CHI SQUARE

Q2. Market line female oxygen $= \frac{296}{449}$ $\hat{p}_1 =$ Get for all

Market line for male and yes ... organic vs non organic

$= \frac{172}{202} \hat{p}_2 =$

$(P_{female} - P_{male}) \pm 1.98 * SE(P_{female} - P_{male})$  $df = 651-2$

$\left(\frac{296}{449} - \frac{172}{202}\right) \pm 1.96$  $\hat{p} \sqrt{\frac{1}{n}+1}$  $t_{critical\ value}$

$SE = \sqrt{\frac{\hat{p}_1 * (1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2 * (1-\hat{p}_2)}{n_2}}$

$\sqrt{\frac{\frac{296}{449}\left(1-\frac{296}{449}\right)}{449} + \frac{\frac{172}{202}\left(1-\frac{172}{202}\right)}{202}}$

$= 0.03356$

$= 0.034$

4. $-0.192 \pm 1.96(0.034)$

$-0.192 \pm 0.0666$

$-0.2586$ btw   $-0.1258$

the negative end of this interval means that market line
for male > market line for female buying organic

2b H₀ There is no association between how voted (buying
organic or not) and gender (male or female) in the population

H₁ There is an association between the two variables (whether you
buy organic) and gender (male or female) in the population

c expected value $= \frac{(Sum\ of\ row)^2 (Sum\ of\ column)}{Total\ population}$

It is the weighted average of all possible value that random
variable can take

2d  Chi square = 25.479  df = 1
      result  is  < 0.001
            probability value  0.001 and 0

P value (0.000) is within the interval - means
we accept the null hypothesis

3  Two Tailed T-Test.
A - Each dot represents a specific number of observations
- The dots are stacked in a column over a category
  the height of the column represents the
  absolute frequency of observation in the category
- X axis represent length
- Y axis represent sample 1 and two

B  Stand error of the means measures the variability of the sample
   = 1√T  Helps determine the difference between more than
   one  sample  of information  standard error is a term that measures
        SE = σ/√n  the accuracy with which a sample represents a population. In stat
                a sample mean deviates from the actual mean of a population,
                this deviation is the SE.

C  H₀  μ(0) - μ(1) = 0  (no difference in means)
   H₁  μ(0) - μ(1) ≠ 0  (difference in means).

   95% Confidence interval for μ(0) - μ(1) = (1.4260 and 2.2050)

   P value is less than 0.005  we can conclude we reject H₀

   T test is measured by ratio  difference between two mean
                                measure of variability or dispersion of data
   - T value is positive, first mean is larger

3c look up t-table with $\alpha = 0.05$ and $df = 149$

$1.97 < t < 1.97$

T (28.41) is outside confidence interval, evidence against Ha

P-value of

p of 0.00 is less than 0.05 and outside the confidence interval, evidence against Ha

3d It is a method of selecting sample members from a large population according to a random starting point and a fixed period interval.

$\frac{1136}{35}$  take every 32$^{th}$ person

4 a Scatterplot:
 - Consists of an x-axis and Y axis, and a series of dots.
 - Each dot represents one observation from a data set.
 - Useful for usually determining the correlation between two variables
 - Independent variable on x-axis (temperature)
 - Dependent variable yield on y-axis

b The equation for yield in grams is $17.0 + 2.0$ (temp c)
 - There is a 2.0 increase for each increase in temperature
 - Coef of constant gives the exact figure for the coefficient in the equation which is rounded to 17 from 17.002
 - Coef of temperature is also rounded up to 2.00 from 1.99517.
 - 17.0 is yield at temp (0).

 - P-value of constant test hypotheses:
   Ho Population Slope = 0        $p = 0.00$ which is $< 0.001$
   H$_1$ Population slope $\neq 0$    Sufficient evidence to reject Ho

- p value in the constant left:

Ho popum intercept = 0

H₁ pop intercept ≠ 0

$p < 0.001$

reject Ho

- S = 4.01967
- S is estimate of standard deviation of Y for fixed x. (SD of residual)

- Residual should be normally distributed
- There should be no relationship between residual and predicted value

- Residual is value of observed cost (actual cost) - predicted value from eqn.
- $R^2$ measures the fit of the model to the data
  98.4% of variance of yield is accounted for

C. T = 75      17 + 2(75)
              17 + 150 = 167 gram