

STAT 2012

1 a ^{Categorical / Nominal} - no order of quantity involved - gender / day of week

Quantitative - quantity involved - underlying scale height / salary

b Graph of length plotted against frequency
 - bars all the same width
 - Shows a) the distribution of variable length

c Mean is the sum of the values of the variables divided by the amount of variables - it is a measure of centre or average

Median is the value in the direct middle when the values are ordered - also a measure of centre

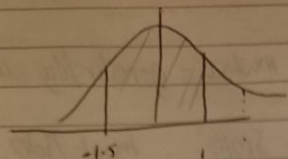
Expected them to be the same if the data is normally distributed implying that they are of similar scale

d. Standard deviation - shows how much variation exists from the average (mean) or expected value low standard deviation indicates that data points tend to be very close to the mean
 It is the positive square root of the variance

Percentile is the value of a variable below which a certain percent of observations fall eg. 20th percentile of the value below which 20 percent of observations can be found

e

mean = 141	standard dev = 2	$1 - 0.93 = 0.07$
$141 - 138 = 3$	$-\frac{3}{2} < z < \frac{3}{2}$	$1 - 0.84 = 0.16$
$143 - 141 = 2$	$1.5 < z < 1$	=



$$P(X < 1) - P(X < -1.5)$$

$$0.84$$

$$P(X < 1)$$

$$1 - P(X < -1.5)$$

$$1 - 0.93 = 0.07$$

$$0.84 - 0.07 = 0.77 = \text{probability}$$

$$2a \quad 1 \pm 1.96 \frac{SD}{\sqrt{651}}$$

$$\frac{183}{651} \text{ purchased organic}$$

$$p = \frac{183}{651} \quad SE(p) = \sqrt{\frac{\frac{183}{651} (1 - \frac{183}{651})}{651}} = 0.017618926$$

$$\frac{183}{651} \pm 1.96 (0.0176)$$

$$\pm 0.0345329$$

$(0.25, 0.32)$ between 25% and 32% buy organic

- b Card type 1 associated with buying organic.
Card type 1 not associated with buying organic

- c Expected value - the weighted average of all possible values that the random variable can take on

$$(\frac{1}{2})^n \quad \% \text{ fine silver occurs } \times \quad \% \text{ fine gold occurs}$$

$$\frac{246}{651}$$

$$\frac{783}{651}$$

fine silver occurs multiplied \times fine gold occurs

$$d \quad Df = 3 \quad \text{chi-square test} = 11.834$$

between 5% and 2.5%

$$p\text{-value} = 0.008$$

chance of accuracy

do not reject which 1) null hypothesis of 0.05

do not reject null hypothesis do not reject with 0.025 and 0.001

4.
If $p < 0.05$ then confidence level we accept is
 $df = 198$
 $1.97 \leq t \leq 1.97$
 t is outside of 1.97 so evidence against H_0
 two long term averages are different

- $p = 0.0001$ which is $p < 0.05$
 - Evidence against H_0

- 95% Interval is 0.5645 and 0.3815 suggest

3c. I would divide 347 by 16 get 21 I would run
 one every 21st minute.

- Scatter plot

4a. Independent variable storage time on x-axis
 - Dependent variable moisture content on y-axis
 - Predict the moisture content by storage time
 - Used to find correlation between two variables

b. Equation to find the moisture at a certain time x
 $2.82 + 0.045 \text{ Storage time}$

- Coef 2.820 and estimate for first coefficient in model eqn
 - Coef 0.045 " " " "
 - 2.820 moisture with 0 hours storage
 - 0.045 - Increase in y per unit increase in x
 - p value on storage time tests the hypothesis:

H_0 Population Slope = 0
 H_1 Population Slope $\neq 0$ $p = < 0.001$ reject H_0

- p value on moisture tests:

H_0 Population Intercept = 0 $p = < 0.001$
 H_1 Population Intercept $\neq 0$ reject H_0

STATS 2012

Q 3 a the box was drawn by graphing one variable against another (graphing 'no' against 'right'). The figure representing 'no' and 'right' were gathered and then graphed

b 95% confidence interval - 95% of the data

$$\mu(\text{no}) - \mu(\text{yes})$$

$$\bar{u}(\text{no}) - \bar{u}(\text{yes})$$

$$124.720 - 130.193 = -0.4730$$

We can conclude that between (-0.5645 and

-0.3815) the mean of (no) is

is between (-0.5645, -0.3815) smaller than mean of (yes)

The t-test assessed whether the means of two groups are different.

- measured by ratio: difference between two means / measure of variability or dispersion of scores

$$\rightarrow \frac{\bar{x}_1 - \bar{x}_2}{SE(\bar{x}_1 - \bar{x}_2)} \text{ (standard error)}$$

$$\rightarrow SE(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{Var_1}{n_1} + \frac{Var_2}{n_2}} \quad \begin{array}{l} Var = \text{Variance} \\ n = \text{Sample Size} \end{array}$$

$$\rightarrow \text{Variance} = (\text{standard deviation})^2$$

$$\rightarrow \text{Final formula} \quad \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{Var_1}{n_1} + \frac{Var_2}{n_2}}} \quad \begin{array}{l} = -0.4730 \\ \sqrt{\frac{(0.5280)^2}{100} + \frac{(0.3280)^2}{100}} \end{array}$$

- If t-value is positive \rightarrow first mean larger, vice versa

- Check alpha level (0.05) and DF = $N_1 + N_2 - 2$

- Look up t table row.

$$\begin{array}{l} -0.4730 \\ 0.04639 \\ = -10.190 \\ \approx -10.20 \end{array}$$

- P value < 0.001 , < 0.5 Accept null hypothesis
long term means are not same

STATS 2012

④ $\sigma^2 = 0.013$ - σ is the standard deviation of the residual

- Residuals should be normally distributed

- There should be no relationship between the residual and the predicted value

- Residual is the value of the observed cost (actual value) - predicted value from equation

- R^2 measures the fit of the model to the data

- 97.8% of the variance of cost is accounted for by the model

c. $2.82 + 0.045t$

$$t = 25 \Rightarrow 2.82 + 0.045(25) = 3.945$$

Company told to the store put the computer in the bin with the data

$$t = 50 \Rightarrow 2.82 + 0.045(50) = 5.07$$

- The t value is outside of the storage time given but is an accurate result.- Contrary to most people's opinion it is not double the value for $t = 25$.