S1 Revision Notes

Representation and summary of Data - Location

Key Words: Continuous, discrete, mean, median, mode, cumulative frequency tables, class width/midpoint, coding

- \blacktriangleright Mean: $\mu = \frac{\sum x}{n} or \frac{\sum fx}{\sum f}$
- ≻ For large data: use coding to work out mean of coded data, then equate to coding for original mean
- Linear interpolation (for some x): using your bivariate date, sketch x against f(x). Plot some a where a < c۶ x and some b where b > x and consider the gradient of $\frac{f(b)-f(a)}{b-a} = \frac{f(x)-f(a)}{x-a}$
 - Used to calculate values in cumulative frequency tables
- Combined mean of set A (size: n, mean: x) and of set B (size: m, mean: y) is given by $\frac{nx+my}{n+m}$ ۶

Representation and summary of Data - Measures of dispersion

iat novia a Gandard deviation Key Words: range, lower/upper quartiles, interquartile range, percentiles, de

- > Variance $\sigma^2 = \frac{\sum x^2}{n} \mu^2 = \frac{\sum fx^2}{\sum f} \left(\frac{\sum fx}{\sum f}\right)^2$
- ≻
- n C ∑f (∑f)
 Standard deviation: how much the member of a scorp differ from the grean value for the group
 For large data: use coding to proceed to the standard deviation of coder data
 Adding result in the DOESNT affect for manipular resilvating does! (Same to the SD!)
 Barcanto constitute to the standard standard for the standard standard to the standard to the standard standard to the standard t ۶
- Percent e spits to data up into xth Percentile $\frac{n}{100}$ Graphiced to calculate percentile range) ۶

Representation of Data

 \geq

 \geq

Key Words: stem and leaf diagram, outliers, box plot, histograms, positive/negative skew,

- \geq An outlier (typically) is any value:
 - Greater than the upper quartile + 1.5 x interquartile range
 - Less than the lower quartile 1.5 x interquartile range 0
 - Box plot: (outlier), (lowest value), (lower quartile), (median), (upper quartile), (highest value), (outlier)
 - Histograms (area is proportional to frequency):
 - \therefore frequency density \times class width = frequency
 - Frequency density always on the y-axis
- \geq Skew is calculated as follows

$$\circ \qquad Q_2 - Q_1 < Q_3 - Q_2 = +ve$$

- $Q_2 Q_1 > Q_3 Q_2 = -ve$
- Mode < mean = +ve \circ

$$\circ \quad Mode > mean = -ve$$

- Using $x = \frac{3(mean-median)}{standard deviation}$ if
 - x = 1 = perfect positive skew
 - x = -1 = perfect negative skew



Probability

Key Words: Venn diagrams, Mutually Exclusive, Exhaustive, Independent, Complementary, Conditional Probability

- \geq Mutually Exclusive: events that cannot happen at the same time e.g. passing and failing an exam
 - \circ $P(A \text{ or } B) = P(A \cup B) = P(A) + P(B) P(A \cap B)$ if NOT ME
 - $\circ \quad P(A \text{ or } B) = P(A \cup B) = P(A) + P(B) \text{ IF ME} (P(A \cap B) = 0)$
- \triangleright Exhaustive: events where all possible outcomes are included e.g. throwing a head or a tail on a fair coin $\circ P(A \text{ or } B) = 1$
- Independent: one event has no effect on another event occurring e.g. throwing a 1 or a 2 on a fair dice \geq

- \circ $P(A \text{ and } B) = P(A \cap B) = P(A) \times P(B)$ if independent events
- *Complementary:* P(A') = 1 P(A)
- Conditional Probability: $P(A \ given B) = P(A|B) = \frac{P(A \cap B)}{P(B)}$
- Multiplication rule: $P(A \text{ and } B) = P(A \cap B) = P(A) \times P(B|A) = P(B) \times P(A|B)$ \geq
- Addition rule: $P(A \cup B) = P(A) + P(B) P(A \cap B)$ \geq

Correlation_ Works tter diaaram. PMCC "r". codina

$$r = \frac{s_{xy}}{|s_{xxxyy}|} \quad s_{xy} = \sum xy - \frac{\sum x \sum y}{n} \quad s_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} \quad s_{yy} = \sum y^2 - \frac{(\sum y)^2}{n}$$

- -1 < r < 1 where:
 - -1 is perfect negative **linear** correlation, 0 is no linear correlation and 1 is perfect positive 0
- Use coding to deal with easier values! The PMCC is not affected by any coding

Regression

Key Words: independent (explanatory) variable, dependent (response) variable, residuals, regression line, interpolation, extrapolation

- A distance, e, from a point (on scatter diagram) to the of best fit is called a residual
 - Outliers have relatively large residuals 0
 - The line of best fit aims to minimise $\sum e^2$ called the regression line of y on x 0
 - Equation of the regression line: y = a + bx, where $b = \frac{s_{xy}}{s_{xx}}$ and $a = \bar{y} + b\bar{x}$ 0
- Interpolation: estimate values within the range of data
- Extrapolation: estimate values outside the range of data

Discrete Random Variables

Key Words: random variable, discrete random variable, cumulative frequency distribution

- For discrete random variable: $\sum P(X = x) = 1$ \geq
- ۶ In a cumulative frequency distribution $F(x) = P(X \le x)$
- $E(X) = \sum x P(X = x)$
- $Var(X) = E(X^{2}) (E(X))^{2}$
- E(aX + b) = aE(X) + b
- $Var(aX + b) = a^2 Var(X)$
- For a Discrete uniform distribution over the values 1,2,3,..,n:

•
$$E(X) = \frac{n+1}{2}$$

• $Var(aX + b) = \frac{(n+1)(n-1)}{12}$

Normal Distribution

Key Words: random variable, mean, standard deviation

 \blacktriangleright If $X \sim N(\mu, \sigma^2)$

•
$$Z \sim N(0,1^2)$$
 where $Z \sim \frac{x-\mu}{\sigma}$ (called standardising)

$$\blacktriangleright \quad \Phi(z) = P(Z < z)$$

 \geq $P(X \le Z \le Y) = \Phi(Y) - \Phi(X)$