

Data Collection

Keywords

Census: measures every member of a population

+ accurate result - expensive / testing may destroy

Sampling units: individuals of a population

Sampling frame: list of sampling units

Random Sampling

i) Simple Random Sampling

? same chance of being selected

use random number / lottery sampling

+ bias free - need sampling frame

ii) Systematic Sampling

? take every k^{th} unit, $k = \frac{\text{pop}}{\text{sam}}$

pick random number between 1 and k for starting point

+ quick to use - need sampling frame

iii) Stratified Sampling

? sample represents groups (strata) of population

$\frac{\text{sam}}{\text{pop}} \times \text{strata}$ for each strata, pick randomly

+ reflects population - pop. must be classified in strata

Non Random Sampling

i) Quota Sampling

? like stratified, but strata filled by interviewer/researcher

+ no sampling frame needed - non-random, potential bias

ii) Opportunity Sampling

? quota filled by those available at the time

+ easy / cheap - unlikely to be representative

Note, there are other advantages and disadvantages

Types of Data qualitative non-numerical

quantitative numerical either discrete or continuous

Large Data Set

UK Stations	When?
Coastal	May-Oct '87 and '95
South	1:0. * Only 6 months
International Stations	
1) Perth, Australia	sun snow
2) Beijing, China	sun clouds snow
3) Jacksonville, FL, USA	sun clouds Oct '87 Oct '95

Data

Rainfall, "tr" means trace, treat as 0.025mm in calculations

n/a means reading not available, so can't be used in sample

Cloud cover oktas, discrete integers 0-8

Max Gust knots, 1kn=1.15mph, Great Storm Oct 15th/16th '87

Location KNOW YOUR CALCULATOR!

Mean

$$\bar{x} = \frac{\sum x}{n} \quad \bar{x} = \frac{\sum fx}{\sum f}$$

Listed Data

$$\text{Position of } Q_1 = \frac{n}{4} \quad \text{Med} = \frac{n}{2} \quad Q_3 = \frac{3n}{4}$$

If a decimal, round up ↑

If whole, find midpoint with next one

Grouped Data

$$\text{Position of } Q_1 = \frac{n}{4} \quad Q_2 = \frac{n}{2} \quad Q_3 = \frac{3n}{4}$$

Percentiles eg. 57th percentile, $P_{57} = 0.57 \times n$

Deciles, 10% chunks, $D_3 = P_{30}$, $0.3 \times n$

Do not round, use linear interpolation

Linear Interpolation

weight, nearest kg	10-12	12.5-15	15-18
frequency	5	8	7
cumulative freq.	5	13	20

$$Q_2 = 12.5 + \frac{5}{8} \times 3 = 14.375 \text{ kg}$$

Spread KNOW YOUR CALCULATOR!

Interquartile Range

$$\text{IQR} = Q_3 - Q_1 + \text{ignores extremes}$$

Interpercentile Range

$$\text{e.g. 10th to 90th IPR} = P_{90} - P_{10}$$

Variance, σ^2 / Standard Deviation, σ

$$\sigma^2 = \frac{\sum x^2}{n} - \bar{x}^2 \quad \text{MSMSM}$$

$$= \frac{\sum fx^2}{\sum f} - \bar{x}^2 \quad \text{Note: } \sum fx^2 \neq (\sum fx)^2$$

"mean of squares minus square of means"

Coding

$$\text{If } y = ax + b$$

$$\bar{y} = a\bar{x} + b \quad \sigma_y = a\sigma_x$$

Representation

Cumulative Frequency / Box Plots



Histograms

continuous data
no gaps

$$\text{frequency density} = \frac{\text{freq}}{\text{class width}}$$

area = freq \times k

Comparisons

i) location ii) spread

Regression + Correlation KNOW YOUR CALCULATOR!

Product Moment Correlation Coefficient

$$\text{PMCC}, r \quad -1 \leq r \leq 1$$

Measures strength and +/- of correlation

Regression Line

$$\text{Line of best fit } y = a + bx$$

a is y when x=0

b is how much y changes with x

✓ interpolation: estimating inside data range

✗ extrapolation: estimating outside data range

Exponential/Non Linear Models

$$\text{If } y = ab^x$$

$$\ln y = \ln a + x \ln b \quad \text{Take logs}$$

$$\text{If } y = ax^n$$

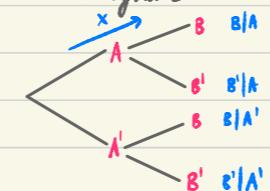
$$\ln y = \ln a + n \ln x$$

Probability KNOW YOUR CALCULATOR!

Venn Diagrams



Tree Diagrams



Mutually Exclusive

If A and B are mutually exclusive

$$P(A \cap B) = 0$$

$$P(A \cup B) = P(A) + P(B)$$

Independence

If A and B are independent

$$P(A \cap B) = P(A) \times P(B)$$

$$P(A|B) = P(A) \quad \text{Can't tell from Venn}$$

Conditional Probability

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$



Reduce the whole space to 1

Addition Law

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Discrete Uniform Distribution KNOW YOUR CALCULATOR!

Probabilities of outcomes all equal

$$\text{e.g. } K \text{ is cloud cover, measured in oktas}$$

x	0	1	2	3	4	5	6	7	8
P(x=x)	1/8	1/8	1/8	1/8	1/8	1/8	1/8	1/8	1/8

Binomial Distribution

$$X \sim B(n, p)$$

$$P(X=r) = \binom{n}{r} p^r (1-p)^{n-r}$$

When to use?

Fixed number of trials, n

Fixed probability of success, p

Independent trials

Two outcomes, success/failure

Cumulative Probabilities

$$P(X < 5) = P(X \leq 4)$$

$$P(X > 3) = 1 - P(X \leq 3)$$

$$P(6 < X \leq 10) = P(X \leq 10) - P(X \leq 6)$$

Inverse Probabilities

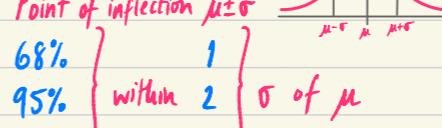
Use tables!

Normal Distribution

Continuous Random Variable, Y

$$Y \sim N(\mu, \sigma^2)$$

Point of inflection $\mu \pm \sigma$



Finding Probabilities and Inverse Normal

$$y. P(Y > 122), P(Y < a) = 0.2$$

Standard Normal Distribution

$$Z \sim N(0, 1)$$

$$\text{Coding, } Z = \frac{Y - \mu}{\sigma}$$

Missing μ , σ or both

Use coding, and sum eq. for both

Approximating Binomial as Normal

If n is large

If $p \approx 0.5$

$$\mu = np \quad \sigma = \sqrt{np(1-p)}$$

Continuity Correction

If approximating binomial as normal

discrete → continuous

$$\text{e.g. } P(X > 5) = P(Y > 5.5)$$

$$P(3 < X \leq 11) = P(3.5 < Y \leq 11.5)$$

Hypothesis Testing

Definitions

Null Hypothesis: H_0 , what we assume to be true