

CH. 1: Data Analysis in an Evidence-Based Practice Environment

Independent variable:

Influence on, an outcome. eg. smoking.

Dependent variable:

The outcome of interest, hypothesized to depend on, or be caused by, the independent variable. eg. Lung cancer, Gender (IV) impact self-esteem (DV).

Level of Measurement:

A system of classification with four types of measurement rules that affect the kind of statistical analysis that is appropriate:

- **Nominal:** Eg. hair colour, gender, religion
- **Ordinal:** Eg. degree of pain, scale of self-esteem
- **Interval:** Eg. temperature
- **Ratio:** Eg. medication dose (number of milligrams, number of pills), pulse

Descriptive statistics:

Describe and summarize data about the sample

Eg. Percent female in the sample, average weight of participants

Inferential Statistics: (what we learn from this course)

- **Parameter 参数:** Eg: Average daily caloric intake of all 10-year-old children in New York
- **Statistic:** Eg: Average daily caloric intake of 300 10-year-old children from three particular NY schools

Inferential statistics, based on laws of probability, help researchers draw objective conclusions about a population, using data from a sample.

Often used to test hypotheses (predictions) about relationships between variables

CH. 2: Frequency Distributions: Tabulating and Displaying Data

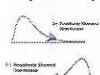
Frequency distribution: Help impose order on the data.

Data Value	Frequency (f)	Percentage (%)	Cumulative Percentage
1	10	10.0	10.0
2	20	20.0	30.0
3	40	40.0	70.0
4	15	15.0	85.0
5	15	15.0	100.0
TOTAL	100	100.0	

10+20
10+20+40
10+20+40+15
10+20+40+15+15

Shapes of Distributions:

- **Modality:** most frequently number
- **Symmetry:**



Positive skew: Longer tail trails off to right (fewer people with high values, like for income)
Negative skew: Longer tail trails off to left (fewer people with low values, like age at death)

Eg. Skewness Index Examples *

Skewness index = 0.80

Standard error = 0.33 → significant skew

Skewness index = -0.72

Standard error = 0.34 → 0.34 x 2 = 0.68 significant skew

CH. 3: Central Tendency, Variability, and Relative Standing

Central Tendency: (specific)

Three alternative indexes:

- The mode:
- The median:
- The mean (N):

Scale	Central Tendency Index	Variability Index
Nominal	Mode	--
Ordinal	Median	Range, IQR
Interval and ratio	Mean	Standard Deviation, Variance

The Mean

Most frequently used measure of central tendency—usually preferred for interval- and ratio-level data.

$M = \frac{\sum X}{N}$

$M = \bar{x}$ (simple mean, x bar)

μ (population mean, mu)

The Range

The difference between the highest and lowest value in the distribution *

110 120 130 140 150 150 160 170 180 190

range = 190 - 110 = 80

The Interquartile Range

Based on quartiles

- Lower quartile (Q₁): Point below which 25% of scores lie — median
- Upper quartile (Q₃): Point below which 75% of scores lie
- Important in evaluating outliers

$IQR = Q_3 - Q_1$

The Standard deviation (SD - Group data):

An index that conveys how much, on average, scores in a distribution vary.

"How far it can go from mean".

$SD = \sqrt{\frac{\sum (X - \bar{X})^2}{N - 1}}$

ex. $SD = \sqrt{\frac{20}{5-1}} = \sqrt{\frac{20}{4}} = \sqrt{5} = 2.236$

X	X - X̄	(X - X̄)²
6	-1	1
7	0	0
7	0	0
7	0	0
7	0	0
7	0	0
8	1	1
9	2	4
9	2	4
10	3	9
10	3	9
11	4	16
11	4	16
12	5	25
12	5	25
13	6	36
13	6	36
14	7	49
14	7	49
15	8	64
15	8	64
16	9	81
16	9	81
17	10	100
17	10	100
18	11	121
18	11	121
19	12	144
19	12	144
20	13	169
20	13	169
21	14	196
21	14	196
22	15	225
22	15	225
23	16	256
23	16	256
24	17	289
24	17	289
25	18	324
25	18	324
26	19	361
26	19	361
27	20	400
27	20	400
28	21	441
28	21	441
29	22	484
29	22	484
30	23	529
30	23	529
31	24	576
31	24	576
32	25	625
32	25	625
33	26	676
33	26	676
34	27	729
34	27	729
35	28	784
35	28	784
36	29	841
36	29	841
37	30	900
37	30	900
38	31	961
38	31	961
39	32	1024
39	32	1024
40	33	1089
40	33	1089
41	34	1156
41	34	1156
42	35	1225
42	35	1225
43	36	1296
43	36	1296
44	37	1369
44	37	1369
45	38	1444
45	38	1444
46	39	1521
46	39	1521
47	40	1600
47	40	1600
48	41	1681
48	41	1681
49	42	1764
49	42	1764
50	43	1849
50	43	1849
51	44	1936
51	44	1936
52	45	2025
52	45	2025
53	46	2116
53	46	2116
54	47	2209
54	47	2209
55	48	2304
55	48	2304
56	49	2401
56	49	2401
57	50	2500
57	50	2500
58	51	2601
58	51	2601
59	52	2704
59	52	2704
60	53	2809
60	53	2809
61	54	2916
61	54	2916
62	55	3025
62	55	3025
63	56	3136
63	56	3136
64	57	3249
64	57	3249
65	58	3364
65	58	3364
66	59	3481
66	59	3481
67	60	3600
67	60	3600
68	61	3721
68	61	3721
69	62	3844
69	62	3844
70	63	3969
70	63	3969
71	64	4096
71	64	4096
72	65	4225
72	65	4225
73	66	4356
73	66	4356
74	67	4489
74	67	4489
75	68	4624
75	68	4624
76	69	4761
76	69	4761
77	70	4900
77	70	4900
78	71	5041
78	71	5041
79	72	5184
79	72	5184
80	73	5329
80	73	5329
81	74	5476
81	74	5476
82	75	5625
82	75	5625
83	76	5776
83	76	5776
84	77	5929
84	77	5929
85	78	6084
85	78	6084
86	79	6241
86	79	6241
87	80	6400
87	80	6400
88	81	6561
88	81	6561
89	82	6724
89	82	6724
90	83	6889
90	83	6889
91	84	7056
91	84	7056
92	85	7225
92	85	7225
93	86	7396
93	86	7396
94	87	7569
94	87	7569
95	88	7744
95	88	7744
96	89	7921
96	89	7921
97	90	8100
97	90	8100
98	91	8281
98	91	8281
99	92	8464
99	92	8464
100	93	8649
100	93	8649

Variance 變動

An important variability concept in inferential statistics, but not used descriptively.*

$\frac{\sum (X - \bar{X})^2}{N - 1}$



Percentiles and Outliers

Outliers are often defined in relation to percentiles

- A **mild outlier** is a score that is between 1.5 and 3.0 times the value of the IQR, below Q₁ or above Q₃
- An **extreme outlier** is a score that is greater than 3.0 times the value of the IQR, below Q₁ or above Q₃

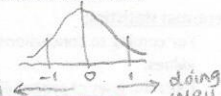
Standard Score (Z score - Individual data)(measure standard deviation)

Another index of relative standing helpful in interpreting raw scores.

A score expressed in standard deviation units, in relative distance from the mean.*

$Z = \frac{X - \bar{X}}{SD}$

ex. $\bar{X} = 500$, $SD = 100$, $X = 600$
 $Z = \frac{600 - 500}{100} = 1$



CH. 4: Bivariate Description: Crosstabulation, Risk Indexes, and Correlation

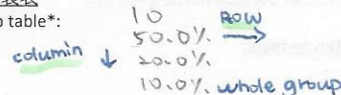
Bivariate descriptive statistics:

Used to describe relationships between two variables.

Eg. Height and weight, Smoking status and lung cancer incidence

Crosstabulation 交叉製表

Example of a Crosstab table*:



Risk Index Scenarios

Eg*:

Independent v.	Cancer (column) dependent v.	
	Yes	No
Smoker	a	b
Control	c	d

ex control	swelling		100
	Yes	No	
16mm	a	b	100
25mm	c	d	100

Absolute Risk: Absolute risk is the proportion of people with a negative outcome.*

$ARE = \frac{a}{a+b}$

$ARE = \frac{20}{20+80} = 0.2 = 20\%$

$AR_{NE} = \frac{c}{c+d}$

$AR_{NE} = \frac{10}{10+90} = 0.1 = 10\%$

Absolute risk reduction: The absolute difference between the two risk groups.*

$ARR = ARE - AR_{NE}$

$ARR = 0.2 - 0.1 = 0.1 = 10\%$

Relative Risk: The ratio of absolute risks (adverse outcomes) in the two groups.*

$RR = ARE \div AR_{NE}$

$RR = \frac{0.2}{0.1} = 2.00$

Relative Risk Reduction: The proportion of baseline risk that is reduced through non exposure (or receipt of an intervention).*

$RRR = ARR \div ARE$

$RRR = \frac{0.1}{0.2} = 0.5$

Odds: The proportion of people in each risk group who have the adverse outcome, relative to the proportion who do not.*

$Odds_E = \frac{a}{b}$

$Odds_E = \frac{20}{80} = 0.25$

$Odds_{NE} = \frac{c}{d}$

$Odds_{NE} = \frac{10}{90} = 0.111$

Odds Ratio: The ratio of the two odds.*

$OR = \frac{Odds_E}{Odds_{NE}}$

$OR = \frac{0.25}{0.111} = 2.25$

Number Needed to Treat: Estimate of how many people would need to avoid the exposure (or get a treatment) to prevent one negative outcome.*

$NNT = \frac{1}{ARR}$

$NNT = \frac{1}{0.1} = 10$

Correlation: Correlations between two quantitative variables can be graphed in a scatterplot.*

maximum number: +1 positive correlation

-1 negative correlation

Correlation Coefficients

- Relationships between two variables.
- Most widely used correlation coefficient: Pearson's product moment correlation coefficient.
- Often called **Pearson's r**.
- Pearson's r is computed with variables that are **interval- or ratio-level** measures.
 - 1.00 = Perfect positive relationship
 - .35 = Weak/moderate positive relationship
 - .00 = No relationship
 - -.20 = Weak negative relationship
 - -.70 = Strong negative relationship

r²

Indicates the proportion of variability in one variable accounted for or explained by the second variable.

Eg. r between height and weight = .60*

$r^2 = .36 \rightarrow 36\%$ the variation in weight is accounted for by height.

64% of variation in weight is accounted for by other factors.

significance: reject (null hypothesis)
 not significance: retain

Correlation Matrix

An efficient way to display several correlation coefficients.

significance correlation
 0.05
 0.04
 0.03
 0.02
 0.01
 0.000001 ← even more significance

scale: interval + ratio
 sig = significance

Pearson correlation: -0.179
 sig: .000
 N: 919
 Positive tick
 number of people in study.

CH. 5: Statistical Inference

Descriptive statistics: For describing samples.

Inferential statistics:

- For coming to conclusions about what is probably true in a population, based on sample values.
- Uses the laws of probability to provide guidance on what is probably true.

Probability of an event (p) is expressed as a proportion

The probability of drawing a red card from a normal shuffled deck:

Ways "red" event can occur = 26
 Total number of possibilities = 52
 $p = .50$: There is a 50-50 chance that the card will be a red suit

Multiplicative law provides this formula:

$p(A \text{ then } B) = p(A) \times p(B)$

The probability of drawing two red cards consecutively:

$p(\text{red, then red}) = p(\text{red}) \times p(\text{red})$
 $p(\text{red, then red}) = .50 \times .50 = .25$
 $p = .25$: There is a one in four chance of drawing two red cards in a row

or (adding)
 or then (multiply)

What is the probability of drawing 10 red cards in a row if the null hypothesis of a fair deck is true?

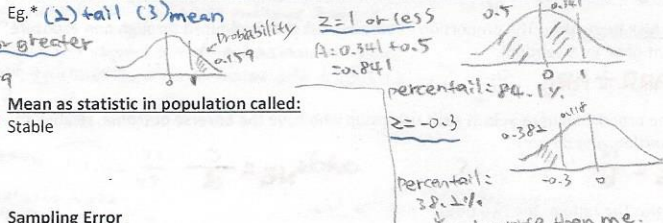
$p = .50 \times .50 \times .50 \times .50 \times .50 \times .50 \times .50 \times .50 \times .50 \times .50 = .001$

Only 1 in 1,000 draws of 10 cards from a fair deck would yield all red cards

Probability distributions

- Similar to frequency polygons (or histograms).
- They graph the probabilities of all events that could occur
 So, the total area of a probability distribution = 1.0

Probability density function = Probability distribution for continuous variables



Mean as statistic in population called:

Stable

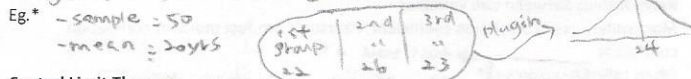
Sampling Error

- Sample means from a population tend to fluctuate from one sample to another because of sampling error.
- Never be exactly 100.0
- Sample mean might be 98.2 or 101.6 or 99.7

Sampling distribution is a type of probability distribution

Sampling distribution of the mean

- A type of probability distribution.
- It's the distribution of an infinite number of sample means from the population.



Central Limit Theorem

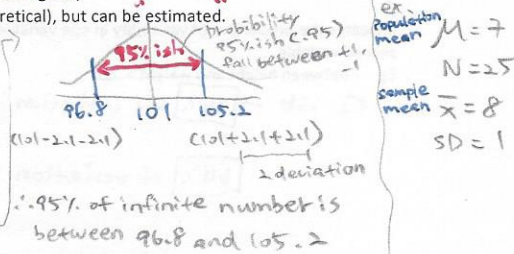
Mathematic formulation, shows that the mean of a sampling distribution of the mean always equals the population mean.

Standard error of the mean (SEM)

- Standard deviation of a theoretical sampling distribution.
- Larger the SEM, the less likely it is that a sample mean is a good estimate of the population mean. (Less: estimate will be good)
- SEMs are never known (due to theoretical), but can be estimated.

Formula for SEM Estimate*

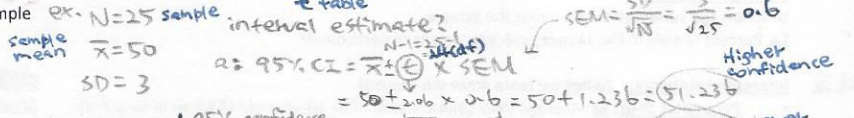
SEM estimate = $\frac{SD}{\sqrt{N}}$
 SD → sample's standard deviation
 N → sample size
 eg. $N = 25$
 $\bar{x} = 101$
 $SD = 10.5$
 $SEM = \frac{10.5}{\sqrt{25}} = 2.1$



Statistical Inference 推理 Approaches

- Hypothesis testing:** Used to for nursing research, now has change.
- Parameter 参数 estimation:** Used to estimate a population value—e.g., a mean, percentage, or odds ratio.
 - Point estimate:**
 - Calculation of a "single value" as the estimate of the parameter.
 - Simply the value of the descriptive statistic, like a mean.
 - Interval estimate:**
 - More accurate! 95% of population fall between these values!
 - Provides a range of values within which the population value has a specified probability of lying.
 - Involves constructing "confidence intervals" around the point estimate.

95% CI = $\bar{x} \pm 1.96 \times SEM$
 SEM = $\frac{SD}{\sqrt{N}}$



t Distribution

- Bell shaped and symmetric.
- A sample size increase, the t distribution is very close to a "normal" distribution.
- Small samples, the tails of a t distribution are "fatter".

CI's around Proportions

- Frequently computed around proportions/percentages and risk indexes like Relative Risk and the Odds Ratio.
- Binomial distribution.
- Almost always done by computer.
- The larger the sample size, the smaller the CI.

Hypothesis Testing

- Common use (retain/reject).
- Second broad approach to statistical inference.
- Help researchers make objective decisions about accepting or rejecting a null hypothesis.
- Eg. Cigarette smoking is unrelated to lung cancer (always start assumed to be innocent!) Cigarette smoking is related to lung cancer

There is always a risk of error (2 types)

- Type I error (*the only type of error can make):**
 - The null hypothesis is really "true" in the population, but the researcher "rejects" it.
 - Incorrect reject, which shouldn't!
 - Controlled through the level of significance (alpha, α) $\alpha = 0.05$, 5% error
- Type II error:**
 - The null hypothesis is really "false" in the population, but the researcher "accepts" it.
 - Incorrect retain, which shouldn't!
 - Probability of committing a Type II error is called beta (β)

$\beta = P(\text{not } H_1)$ Type II error
 power = $1 - \beta$
 bigger sample size = more power!

Power

- 3 factors influence power:
 - Sample size → the easiest one to control.
 - Effect size
 - Alpha level
- Power and the probability of making a Type II error are inversely related.
- ↓ the probability of making a Type II error, ↑ the probability of making a Type I error

determine type II error ← $\text{Power} = 1 - \beta$ → type II error
 ex. Power = 0.8 (good) $\beta = 1 - 0.8 = 0.2$ (20% type II error)

One-Sample t Tests

- A statistical test that tests the null hypothesis that the population mean is a "specific value".
- Tests whether the sample mean differs from a real or hypothesized population mean. (compare "Sample mean" to "Population mean")

