

CHI Look at data - distribution of 1 variable
 units/cases into a characteristic - observation
 what how many individual - variable
 individuals
 categorical variable
 numerical quantitative variable
 distribution
 what value how often
 tools - nature of variables

Graph.

categorical variable - distribution.
 individual → group
 (count) individuals
 (percent) each cate.
 Pie chart (distrib. cate. pie size - count percent)
 Bar graph (Bar - cate. height - count percent)

Quantitative variable - distribution -
 take numerical values (what value how often)
 Histogram (Bar - class height - # individuals)
 Stemplots - original value
 Exam (symmetric, right-skewed, left-skewed)
 outlier - pattern.

Numbers

Center - mean
 outlier
 Median.

Spread - range.
 variance, standard deviation

Quartile ranked data - 4 equal parts
 Q_1, Q_2, Q_3 - Index → value.

12.5%	25%	25%	25%
Q_1	Q_2	Q_3	

Inter-quartile range: middle 50%
 $Q_3 - Q_1$

Measure - Center, spread, outlier.

Median + IQR → skew distribution
 outliers.

mean + standard deviation - (symmetric, outlier).

Numeric summaries

5 Number summary: min, Q_1 (25th percentile), Q_2 (median), Q_3 , max.

plot data - Box plot.

min max
 Q_1, Q_2, Q_3 - IQR
 lower/upper Fence.

spread
 outlier
 Center.
 shape - symmetric.

Ch 2.

Relationships - 2 variables.

Cases / units - info. characteristic.

(what. how many) \rightarrow variable. (cate. quant. values)

(Explanatory (X)

response (Y))

Exam Relationships. \rightarrow cate. variables - contingency table.

$\checkmark \rightarrow$ quant. variables - scatterplot.

Scatterplot.

1 cate. + 1 quant - boxplot side-by-side.

Same individuals \rightarrow each individual - 1 point.

Which variable - which axis: (explanatory - X (label. response - Y (Scale)

Interpret Scatterplot.

(patterns. departures.

direction \times no - nonlinear - linear.
Form \nearrow 70 \searrow 30
no weak moderate strong
Strength 0 - 0.3 - 0.5 - 0.7 - 1
outlier.

+ 1 cate. - diff. plot color / symbol.

numerical Supplement:

correlation coefficient (r) - linear relationship. \checkmark causation. \checkmark cause and effect.

X = working variable \rightarrow Y. | large.

least square regression line - ~~fit~~ \rightarrow fit.

$$(b_1) \quad \bar{x} \quad \bar{y} \quad b_0 = \bar{y} - b_1 \bar{x} \rightarrow \hat{y} = b_0 + b_1 x$$

Coefficient of determination (r^2) \rightarrow fit original data points. (extrapolating interpolating)

overall pattern. - residual $y - \hat{y}$
residual plot.

outlier - influential point.

24.

Quiz Score: 30 out of 30

Random experiment - variable - outcomes - Event

Event - Events

Count
probability

$A \cup B$
 $A \cap B$
 $A' \cap B'$
 $A \cap B = \emptyset$

Variable outcomes count probability

1 single variable

2 cate. variable - outcomes

contingency table

$A \cap B$ - counts - probability
 $A \cap B$

Joint
marginal > contingency table

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

Joint probability $P(A \cap B)$ independent events $P(A \cap B) = P(A) \cdot P(B)$

Marginal probability $P(A)$
 $P(B)$

Multiplication Rule $P(A \cap B) = P(A|B) \cdot P(B)$
 $= P(B|A) \cdot P(A)$

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

trials n	outcomes			Prob. calc.	Prob.
	1	2	3		

independent
dependent

Conditional Probability
Marginal Probability

tree diagram

notation.
事件, 问题
tree diagram. number $P(A)$ $P(B)$
Answer - Step.

Total Probability Rule $P(B) = P(B|A) + P(B|\bar{A})$
Bayes' Theorem
 $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$

Discrete Random Variable

experiment - outcomes - Sample space.
events - variable - take on values.

variable - values - probability.
(list - distribution)
Probability mass function (PMF) $f(x)$
 $f(x) \geq 0$
 $\sum f(x) = 1$
 $f(x) = P(X=x)$

Discrete Random variable
counted
(PMF) $f(x)$

Cumulative Distribution function (CDF) $F(x)$
 $F(x) = P(X \leq x) = \sum_{x_i \leq x} f(x_i)$
PMF - CDF - change.
 $F(x) - P(X \leq x) = 1 - F(x)$

Summary numbers
mean = center / median $\mu = \sum x \cdot f(x)$
variance / standard deviation - dispersion $\sigma^2 = \sum x^2 \cdot f(x) - \mu^2$
 $\sigma = \sqrt{\sigma^2}$
 $\sum x \cdot f(x) = \bar{x} = \mu$
 $\sum x^2 \cdot f(x) = \sum x^2$
 $\sum = -\mu^2 = \sigma^2$

Binomial

Experiment: each trial - outcomes
independent
Trials - outcomes - success - Prob
 n

Success P
 n Binomial (n, p)
 x
 1

$\text{BINOMDIST}(\text{number, trials, Prob, cumulative})$
 $\text{BINOMDIST}(x, n, p, \text{True/False})$
 $X \sim \text{Binomial}(n, p)$

Mean: $\mu = E(X) = np$
 $\sigma = \sqrt{np(1-p)}$

$f(x) = P(X=x) = C_n^x p^x (1-p)^{n-x}$
 $P(X=x)$
at least $z - P(X \geq z)$
more than $z - P(X > z)$
at most $z - P(X \leq z)$
fewer than $z - P(X < z)$
no more than $z - P(X \leq z)$
 $P(z \leq X < 7)$

Geometric Distribution

Quiz Score: **33** out of 33

A trial - outcomes - success - Prob
 Independent

trials x - until 1 success - Prob $P(X=x) = (1-p)^{x-1} \cdot p$

Continuous Random Variables

variable x - Probability Density function $f(x)$
 (Normal Density)

More than 5 $P(X > 5) = 1 - P(X \leq 5)$
 $= P(X=1) + P(X=2) + P(X=3) + P(X=4) + P(X=5)$
 Probability distribution - curve
 Normal distribution $X \sim N(\mu, \sigma)$
 68-95-99.7

Standard Normal distribution -

$$Z = \frac{x - \mu}{\sigma} \sim N(0, 1)$$

Cumulative distribution function.
 $\Phi(z) = P(Z \leq x) = F(z)$

L7 (Binomial)

Binomial \rightarrow Normal (11)

$$\begin{aligned} np &\geq 10 \\ nq &\geq 10 \end{aligned} \quad \left. \begin{aligned} \mu &= np \\ \sigma &= \sqrt{npq} \end{aligned} \right\}$$

Continuity correction

$$P(X < x) \rightarrow P(X \leq x - 0.5)$$

$$P(X > x) \rightarrow P(X \geq x + 0.5)$$

$$Z = \frac{x - np}{\sqrt{npq}}$$

Normal \rightarrow Standard normal $+0.5 \quad -0.5$

Sampling distribution of (sample mean) $\bar{x} = \frac{x - \mu}{\frac{\sigma}{\sqrt{n}}}$
 population - distribution = Normal (sample proportion) $\hat{p} = \frac{x - np}{\sqrt{npq}}$ (suppose)

size n Sample

Sample mean $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$

Probability distribution - Sample distribution of Sample mean

(value interval)

(how often 'probability')

$$\mu_{\bar{x}} = \mu$$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Normal

population property

(Population \rightarrow Sample)

mean $\bar{x}_1, \dots, \bar{x}_n \rightarrow \bar{\bar{x}}$
 Sample distribution interval (8)
 \bar{x} how often - Prob

Normal standard Normal

Normal distribution: population

$$P(Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}})$$

sample

$$\bar{x} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$n \geq 30$$

$$v.$$

Success - class cate.
 population - population proportion $p = \frac{\# \text{Success}}{N}$

Sample - Sample proportion = $\frac{\# \text{Success}}{n} = \hat{p}$

Sampling distribution of Sample proportion.

n large - Sample proportions $\hat{p} \sim \text{Normal distribution}$
 A random sample

Success
 Population - success p
 Sample size n
 Success - $\hat{p}_1, \dots, \hat{p}_n \rightarrow \hat{p}$
 Sampling distribution interval (\hat{p})
 Prob $P(\hat{p})$
 Success - p
 if $\hat{p} = p$ $P(\hat{p}) = P(p)$
 $\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$
 $P(\hat{p}) \rightarrow P(p)$

Confidence Interval

Population parameter $\rightarrow \mu$

Sample \rightarrow past studies, studies
 Standard Normal $\hat{p} \rightarrow z$

$z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}}$ - interval probability

point estimate - 1 sample - observations - Sample mean \bar{x} - x how close.
 interval estimate - interval - point estimate \bar{x}

Confidence interval

point estimate \bar{x}
 Critical value $z_{\alpha/2}$ - 2-sided - $z_{\alpha/2}$
 Margin of error $m = (z_{\alpha/2}) \cdot \frac{\sigma}{\sqrt{n}}$
 interval $(\bar{x} - m, \bar{x} + m)$
 Population parameter $\rightarrow \mu$

point estimate, 1 sample - sample proportion \hat{p}
 Critical value
 Margin of error $m = (z_{\alpha/2}) \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ $\hat{p} \hat{p} - 0.5$

$m = z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

Confidence interval $\hat{p} \pm m$

Hypothesis test.

Population mean μ
 Population proportion p
 Hypotheses

Population $(\frac{\mu}{\sigma})$

Sample n

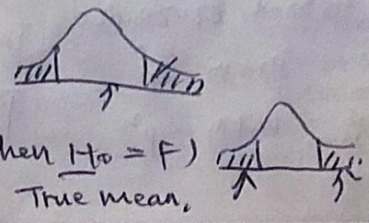
$\bar{x} \rightarrow \hat{p}$

H_0
 H_1 : - 1 or 2 sided \rightarrow left or right tail

- Critical value z_{α} or $z_{\alpha/2}$
- Test statistic: z_0
- Decision \hookrightarrow Rejection region
- P-value

LTP

H_0 T F type I error $\alpha = P(\text{Reject } H_0, \text{ when } H_0 = T)$
 T_0 T type II error $\beta = P(\text{Fail to reject } H_0, \text{ when } H_0 = F)$
 Sample mean True mean



12. 2 sample \rightarrow population proportion p_1, p_2

Group

Sample	variable	Proportion
1	n_1	X_1
2	n_2	X_2

$$\hat{p}_1 = \frac{X_1}{n_1} \rightarrow p_1$$

$$\hat{p}_2 = \frac{X_2}{n_2} \rightarrow p_2$$

Solution 1

Mean $\mu = np = 5(0.7) = 3.5$
 Variance $\sigma^2 = np(1-p) = 5(0.7)(1-0.7) = 1.05$
 Standard Deviation $\sigma = 1.02$

CI:

point estimate

critical value

$\alpha/2$

margin of error

interval

$$\hat{p}_1 \hat{p}_2 \rightarrow \text{CI - 2 sided}$$

$$n = \sum \alpha/2 \cdot \frac{\sqrt{\hat{p}_1(1-\hat{p}_1)}}{n_1} + \frac{\sqrt{\hat{p}_2(1-\hat{p}_2)}}{n_2}$$

Hypothesis test

$$H_0: p_1 = p_2$$

$$H_1: \text{1/2 sided left/right tail}$$

$$\text{critical value } \hat{p}_1 \hat{p}_2 \rightarrow \text{CI}$$

Fig. Z

Test statistics

$$Z_0 = \frac{(\hat{p}_1 - \hat{p}_2) - \mu(\hat{p}_1 - \hat{p}_2)}{\sqrt{\hat{p}_1(1-\hat{p}_1) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad p = \frac{X_1 + X_2}{n_1 + n_2}$$

Decision:

plot + conclusion

P-value

2 sample - population mean μ_1, μ_2

Group

Sample	mean	variance
1	$\bar{X}_1 \rightarrow \mu_1$	σ_1^2
2	$\bar{X}_2 \rightarrow \mu_2$	σ_2^2

CI:

point estimate $\bar{X}_1 - \bar{X}_2$

critical value

margin of error

interval

$$\sigma_1^2 \vee \sigma_2^2 \rightarrow \text{CI - 2 sided}$$

$$n = \sum \alpha/2 \cdot \frac{\sqrt{\sigma_1^2}}{n_1} + \frac{\sqrt{\sigma_2^2}}{n_2}$$

Hypo. Test

Hypo.

H_0

$$H_1: \text{1/2 sided}$$

$$\text{left/right tail}$$

$$\text{critical value } \sigma_1^2 \sigma_2^2 \rightarrow \text{CI}$$

Fig. Z

Test statistics

$$Z_0 = \frac{(\bar{X}_1 - \bar{X}_2) - \mu(\bar{X}_1 - \bar{X}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

critical value $\sigma_1^2 \sigma_2^2 \rightarrow t$

$$df = n_1 + n_2 - 2$$

Test statistics

$$P_0 = \frac{(\bar{X}_1 - \bar{X}_2) - \mu(\bar{X}_1 - \bar{X}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{(\bar{X}_1 - \bar{X}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

critical value $\sigma_1^2, \sigma_2^2 \rightarrow t$

$$df = n_1 + n_2 - 2$$

Test statistics

$$T_0 = \frac{(\bar{X}_1 - \bar{X}_2) - \mu(\bar{X}_1 - \bar{X}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{(\bar{X}_1 - \bar{X}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$\bar{X} - m, \bar{X} + m$$

Critical value: $S_1, S_2 \rightarrow t$ CI - 2 sided - $\alpha/2$

$$m = \sum \alpha/2 \cdot \left(\frac{\sqrt{S_1^2}}{n_1} + \frac{\sqrt{S_2^2}}{n_2} \right)$$

$$\bar{X} - m, \bar{X} + m$$

Critical value $S_1, S_2 \rightarrow t$ CI - 2 sided

$$m = t \cdot \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \quad n_1 = n_2 \quad S_p^2$$

$$= t_{\alpha/2} \cdot S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \quad n_1 \neq n_2 \quad S_p^2$$