

No dumb questions!

Chapter 1 Introduction to Data

1.2.

Observation: Each row in the data set is an observation / case / unit.

Variable: Each column represents some characteristic, varies between observations.

Numerical Variables: the values are numeric, the numbers are meaningful, measuring quantities.

Discrete Numerical:

Continuous Numerical:

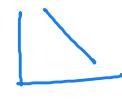
Categorical Variables: the values are categories / types / labels

Ordinal Categorical: the categories are ordered.

Nominal Categorical: " " unordered.

Scatter plot: Use dots to denote data from 2 variables.

Positively associated:  scatterplot shows an upward trend

Negatively associated:  " " downward trend.

Response Variable (Dependent Variable.)

Explanatory Variable (Independent Variable)

1.3

Population: all the units of interest.

Individual Cases / Subjects: one unit in a population.

Anecdotal Evidence: (small sample size. extreme cases)

Sampling "Representative of the population"

{

Simple Random Sampling — "raffle method"

Statified Sampling — "divide and conquer"

Clustering — "mini population"

Multi-stage

1.4 Experiments "Apple to Apple"

1. Controlling
2. Randomization
3. Replication
4. Blocking

Chapter 2

2.1

Scatter Plot — "2 variables"

Dot Plot

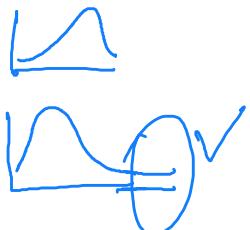
Stacked Dot Plot

Histogram: bin size "same"

} 1 variable

Shape:

{ Symmetric
Left Skewed
Right Skewed



2.2 Center

① "Mean":

$$\text{Sample mean } \bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

Population mean μ

② Median: The number that cuts the data in half. "50th" percentile.

odd # of obs: $\frac{n+1}{2}$ position

even # of obs: $\frac{n}{2}, \frac{n}{2}+1$ take average

③ Mode:

Variation / Spread.

① Variance:

$$\text{Sample Variance } s^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n-1} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Sample Standard deviation $s = \sqrt{s^2}$

Population Variance σ^2

Population St. dev. σ

② IQR = $Q_3 - Q_1$, Interquartile Range.

Robust Statistics : insensitive to outliers.

✓ IQR , Median

X Var, St.dev. mean

2.3 Visualization

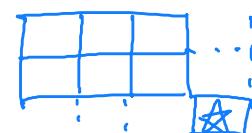
1 Var } Bar Chart
 } Pie Chart

2 Categorical Variables

Contingency table (2×2 , $3 \times 4 \dots$)

Stacked Bar chart

Mosaic Chart



| Categorical Variable + | Numerical Variable.

Side-by-Side Box plot

Stacked Histogram.

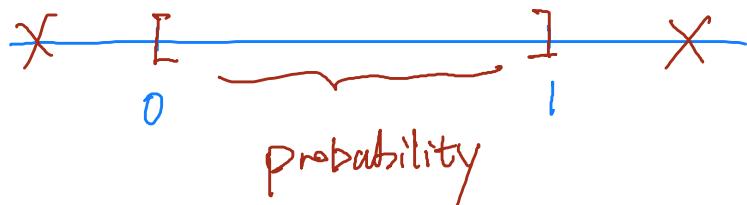
Chapter 3 Probabilities — Part 1.

3.1

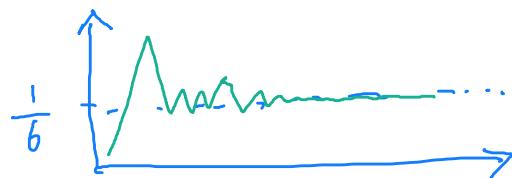
Probability: "infinite number of time" \rightarrow proportion

0 \rightarrow never

1 \rightarrow always



Law of Large Number



Sets and Events

Event: a set outcomes we are interested in

Use A, B, C... $A = \{4, 6\}$ $P(A) = \frac{2}{6} = \frac{1}{3}$

$B = \{2, 4, 6\}$

Disjoint / Mutually Exclusive Events

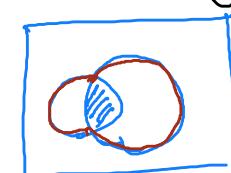
A and B \times disjoint.

General Addition Rule

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

U

Venn Diagram



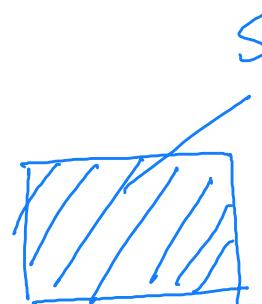
Probability Distribution : a list of all possible results with corresponding probabilities

Price	$< 300K$	$300 \leq$	$< 600K$	$600K - 1000K$	$> 1000K$
-------	----------	------------	----------	----------------	-----------

- ① all results/outcomes are disjoint/mutually exclusive
- ② all possible results/outcomes.
- ③ all probabilities add up to 1.

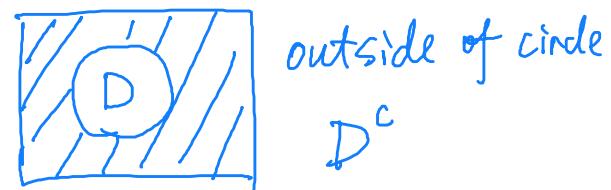
Sample Space : A set of all possible outcomes.

- Denoted by "S"
- $P(S) = 1$



Complement of event D: "not" D ,

- Denoted by $D^c \rightarrow$ not D



- $P(D) + P(D^c) = 1$

- D and D^c are disjoint. $P(D \text{ or } D^c) = 1$

o

Independent events: $P(A)$ does not depend on whether B happens or not.

- $P(A \text{ and } B) = P(A)P(B)$

- If A and B are disjoint, then A and B are NOT indep.

A even number $\{2, 4, 6\}$ 50%

B greater than 3 $\{4, 5, 6\}$ 50%

A - B not disjoint .

$$P(A \text{ and } B) = \frac{2}{6} = \frac{1}{3}$$

$$P(A) \times P(B) = \frac{1}{2} \times \frac{4}{6} = \frac{1}{3}$$

3.2 Conditional Probabilities.

	Bad	Good
Algorithm	Bad	22 X
	Good	1491
		1513

1. ✓ $P(\text{Algorithm say "Good"} | \text{Good}) = \frac{1491}{1513} = 98.5\% \leftarrow$

2. ✓ $P(\text{Algorithm say "Bad"} | \text{Bad}) = \frac{197}{197+112} = 63.7\% \leftarrow 100\%$

$$P(\text{error rate}) = \frac{112+22}{197+112+22+1491} = 7.4\%$$

3. $P(\text{correct rate}) = 1 - P(\text{error rate}) = 93.6\% \leftarrow$

Condition : information we know to be true. (B)

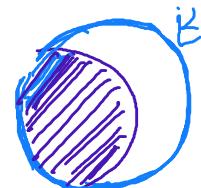
Outcome of Interest: what we want to know. (A)

$$P(\text{outcome of interest} | \text{condition}) = P(\text{Label good} | \text{good})$$

given

Label
Good
A

$$= \frac{P(\text{Label as good})}{P(\text{good videos})} = \frac{\text{Shaded area}}{\text{Total area}}$$



$P(A \cup B)$

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

' U
or

General Multiplication Rule.

$$P(A \cap B) = P(B) \cdot P(A|B) \checkmark$$

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

$$P(A \cap B \cap C) = P(A)P(B|A) \cdot P(C|A \cdot B)$$

Independent

$$P(A \cap B) = \underline{P(A)P(B|A)}^{\text{indep}} = P(A)P(B)$$

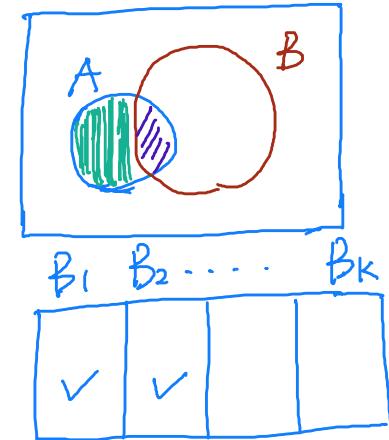
★ $P(B) = P(B|A) = P(B|A^c)$

$$P(A) = \text{shaded area} + \text{unshaded area}$$

$$= P(B)P(A|B) + P(B^c)P(A|B^c)$$



$$= P(B_1)P(A|B_1) + P(B_2)P(A|B_2) + \dots + P(B_k)P(A|B_k)$$



Chapter 3 Part 2

Random Variable:

A variable is a random variable if the value it takes on is a random event. (win/lose election, # of heads if flipping 5 coins).

Expected Value of a R.V. $\rightarrow E(X) \cdot \mu$, mean

X, x_1 with $P(X=x_1)$, x_2 with $P(X=x_2) \dots x_n$ with $P(X=x_n)$

$$E(X)=\mu = x_1 \cdot P(X=x_1) + x_2 P(X=x_2) + \dots + x_n P(X=x_n)$$
$$= \sum_{i=1}^n x_i P(X=x_i)$$

Variance of a R.V.

$$\text{Var}(Y)=\sigma^2 = (x_1-\mu)^2 P(X=x_1) + \dots + (x_n-\mu)^2 P(X=x_n)$$
$$= \sum_{i=1}^n (x_i-\mu)^2 P(X=x_i)$$

$$\text{std}(Y)=\sigma = \sqrt{\text{Var}(Y)}$$

Linear Transformation of a R.V. X

$$Y = aX + b$$

$$E(Y) = aE(X) + b$$

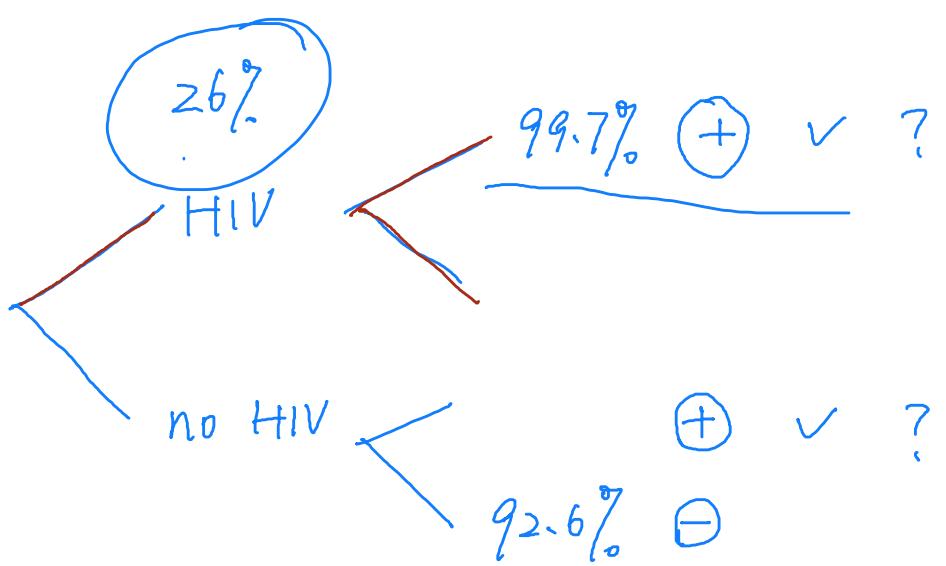
$$\text{Var}(Y) = a^2 \text{Var}(X)$$

$$Z = X_1 + X_2 + X_3$$

$$E(Z) = E(X_1) + E(X_2) + E(X_3)$$

$$\text{Var}(Z) = \text{Var}(X_1) + \text{Var}(X_2) + \text{Var}(X_3)$$

independent X_1, X_2, X_3



(a) $\underline{(26\%)(99.7\%)} + (1 - 26\%)(1 - 92.6\%) = P(+)$

(b) $\frac{(26\%)(99.7\%)}{P(+)} = P(HIV | +)$

Chapter 4 Distributions

1. Bernoulli Distribution
2. Binomial Distribution ←
3. Normal Distribution
1. Bernoulli Distribution

• 30% like blue cheese.

X	1 (Yes)	0 (No)
prob	0.3	1-0.3

$$p = P(\text{success}) = 0.3$$

X	1	0
prob	p	$1-p$

$$E(X) = p$$

$$\text{Var}(X) = p(1-p)$$

2. Binomial Distribution:

3 randomly selected student, X : # of students like blue cheese.

X	0	1	2	3
prob				

$$X = X_1 + X_2 + X_3$$

$$P(X=0) = (1-0.3)(1-0.3)(1-0.3) = (1-0.3)^3$$

Andy \star Brad \star Cathy \star
 ✓ x x

$$P(X=k) = \underbrace{(0.3)(0.3)(0.3)}_k = 0.3^3$$

$$P(X=1) = \frac{1}{3} \cdot (1-0.3)^2 \cdot (0.3)^1$$

$$P(X=1) = \frac{(0.3)(1-0.3)(1-0.3)}{3} \times 3 = (0.3)(1-0.3)^2$$

$$P(X=2) = (0.3)(0.3)(1-0.3) \quad \times 3 = (0.3)^2(1-0.3)^1$$

$$\begin{aligned}
 E(X) &= \sum x_i P(X=x_i) = 0(1-0.3)^3 + 3(0.3)^3 + 1(0.3)(1-0.3)^2 \cdot 3 + 2(0.3)^2(1-0.3) \\
 &= E(X_1 + X_2 + X_3) \\
 &= E(X_1) + E(X_2) + E(X_3) \\
 &= (0.3) + (0.3) + (0.3) \\
 &= (0.3) \times 3
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(X) &= \text{Var}(X_1 + X_2 + X_3) = \text{Var}(X_1) + \text{Var}(X_2) + \text{Var}(X_3) \\
 &= (0.3)(1-0.3) + 0.3(1-0.3) + 0.3(1-0.3) \\
 &= 3(0.3)(1-0.3)
 \end{aligned}$$

Binomial Setting

1. X is # of successes from "n" independent trials. (\checkmark/X)
2. "independent"
3. Probability of success in each trial is constant "p"
4. "n" is fixed.

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

\downarrow
#success

$$E(X) = E(X_1 + X_2 + \dots + X_n) = p + p + p + \dots + p = np$$

$$\text{Var}(X) = \text{Var}(X_1 + X_2 + \dots + X_n) = np(1-p)$$

n

Counting

$$\bullet \binom{n}{k} n \text{ choose } k \quad \frac{n(n-1)\cdots(n-k+1)}{k!} = \frac{n!}{k!(n-k)!}$$

$k! = k(k-1)\cdots 1$

Combination C_n^k

rank

- order k items $K(K-1)(K-2)\dots 1 = k!$
- order k items from n $100 \times 99 \times 98 \dots P_{100}^3$

Permutation. P_n^k

$$3! = 6.$$

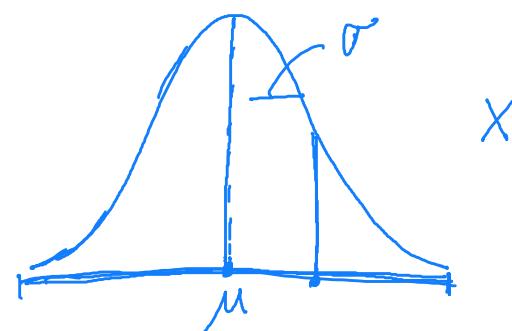
$$\binom{5}{2} = \frac{5 \times 4}{2 \times 1} = 10$$

3. Normal Distribution (Continuous R.U.)

$$\textcircled{x} \quad f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\left(\frac{(x-\mu)}{\sigma}\right)^2}$$

$$X \sim N(\mu, \sigma)$$

Whole area under the curve is 100%.

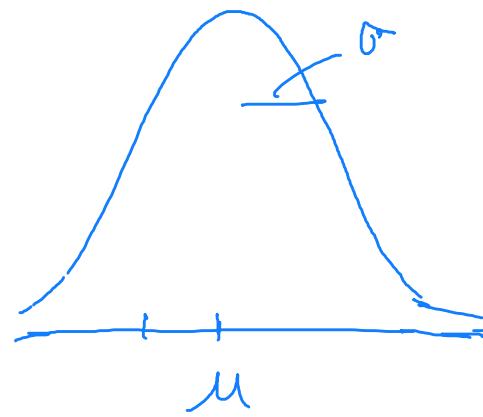


$$P(X \geq \mu) = 50\%$$

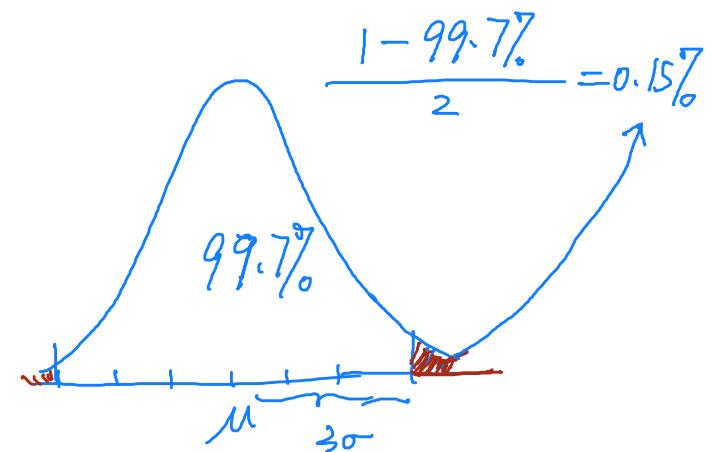
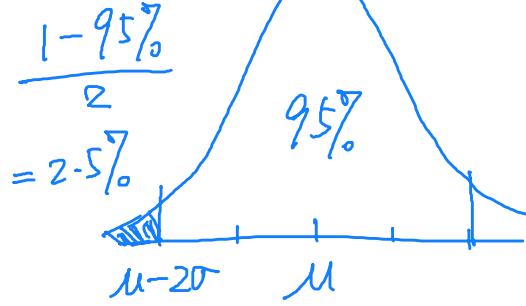
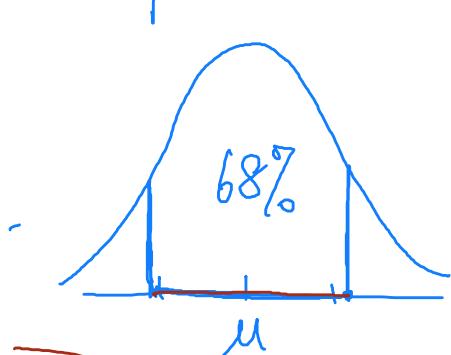
$$P(X = \mu) = 0$$

\uparrow
z-score = how many st.dev is X away from the center.

$$Z = \frac{X - \mu}{\sigma}$$



Empirical Rule $68 - 95 - 99.7\%$



$\boxed{\mu \pm \sigma}$: 1 σ within μ

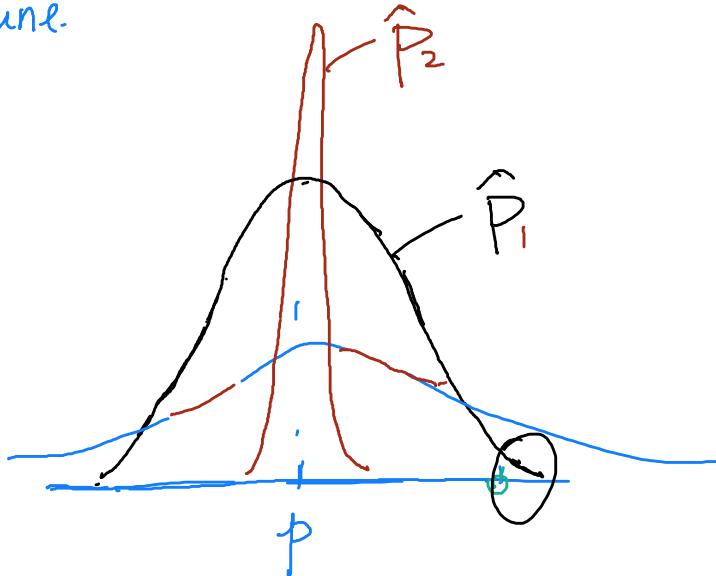
$\mu \pm 2\sigma$: 2 σ within μ

$\mu \pm 3\sigma$

Chapter 5

Trump Approval Rate in 2020 June.

- ① $n=5000 \quad \hat{P}_1 = 41\%$
- ② $n=5000 \quad \hat{P}_2 = 42\%$
- ③ $n=5000 \quad \hat{P}_3 = 38\%$
- ⋮
- $\hat{P}_{1000} =$



• Central Limit Theorem (CLT)

The sample proportion \hat{P} will approximately follow a normal distribution with p as the mean and $\sqrt{\frac{p(1-p)}{n}}$ as the standard deviation if $np \geq 10$ and $n(1-p) \geq 10 \rightarrow$ success/failure condition.

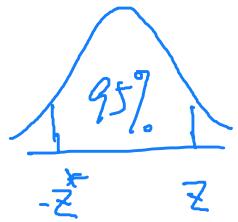
$$\hat{P} \sim N(p, \sqrt{\frac{p(1-p)}{n}}) \text{ if } np \geq 10, n(1-p) \geq 10$$

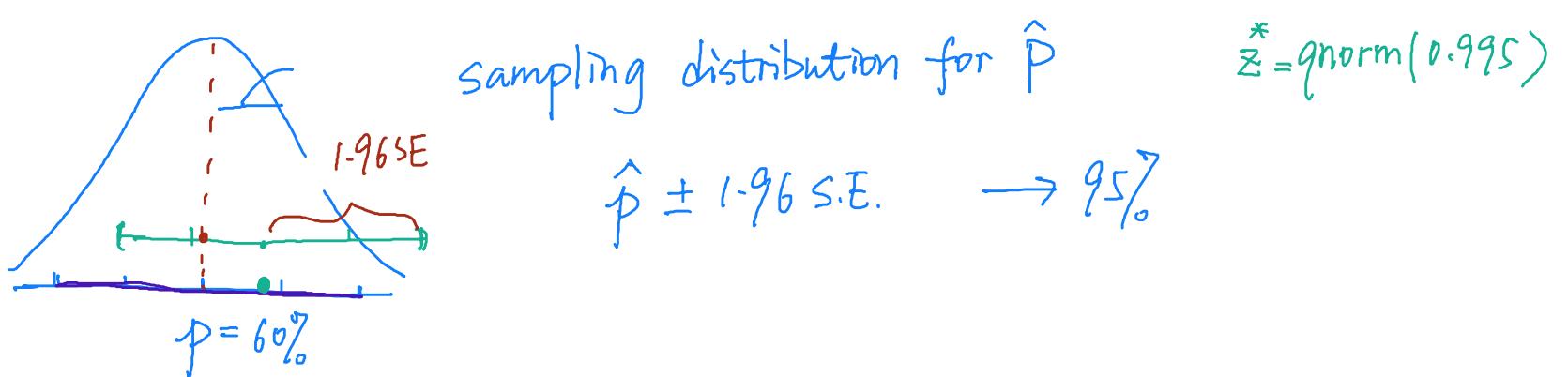
• Confidence Interval

$$\hat{P} \pm \text{margin of error}$$

$$\hat{P} \pm z^* \cdot SE$$

95%	$\rightarrow 1.96$
90%	$\rightarrow 1.645$
99%	$\rightarrow 2.58$



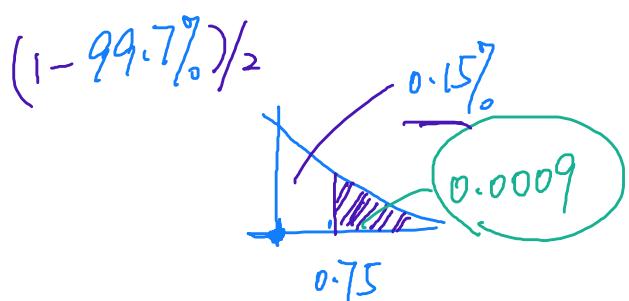


- ① Large $n \rightarrow$ smaller m.o.e.
- ② lower confidence level \rightarrow smaller m.o.e.

Hypothesis Testing

$\left\{ \begin{array}{l} 50\% \text{ — overall death rate.} \\ \underline{40 \text{ patients}} \text{ — 30 died.} \rightarrow 75\% \end{array} \right.$

$\hat{P} \sim N(0.5, \sqrt{\frac{0.5(1-0.5)}{40}} = 0.08)$
 $z = \frac{0.75 - 0.5}{0.08} = 3.125 \leftarrow \text{pnorm}(-3.125) = 0.0009 \times$



$H_0: p = 0.5 \star$ Null hypothesis

$H_a: p > 0.5$ Alternative hypothesis

		H_0	Reject H_0
		true	false.
✓ Reject H_0	true	X I	✓
	false.		X II
Fail to reject H_0	true	✓	
	false.		X II

Type I : H_0 is true, reject H_0 .

Type II : H_0 is false, fail to reject H_0 .

Chapter 6 Inference for Population Proportion

- Hypothesis Test for one "p"
- Hypothesis Test for two "p's" $p_1 - p_2$

Hypothesis Test for one "p"

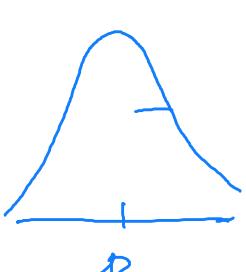
① $H_0: p = p_0$ \times

$H_a: p \neq p_0$ — two sided test

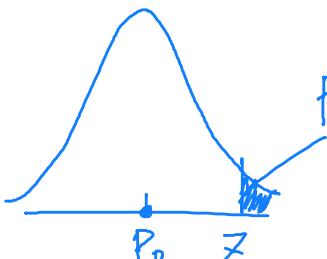
$$\begin{array}{l} p > p_0 \\ p < p_0 \end{array} \quad \left. \begin{array}{l} \text{two sided test} \\ \text{one sided test} \end{array} \right\}$$

② Check CLT condition

- independent sample
- $np_0 \geq 10$ and $n(1-p_0) \geq 10$

③  $\hat{p} \sim N(p_0, \sqrt{\frac{p_0(1-p_0)}{n}})$

$$Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \quad \text{test statistic}$$

④  p-value : how likely to observe some data equally or more extreme than \hat{p} .



⑤ Conclusion: compare p-value with significance level (α) Type I error.

p-value $\geq \alpha \rightarrow$ fail to reject H_0

p-value $< \alpha \rightarrow$ reject H_0

 0.01 0.05

HT for $p_1 - p_2$

① $H_0: p_1 = p_2 \leftarrow$

$H_a: p_2 - p_1 > 0$ trt 2 is better?

trt 1: $n_1 = 109$ $\hat{p}_1 = 60\%$

trt 2: $n_2 = 91$ $\hat{p}_2 = 65\%$

② Check CLT for \hat{P}_1, \hat{P}_2

independence: within each group, between each group

success/failure: $n_1 \hat{p}^* \geq 10$ $n_1(1-\hat{p}^*) \geq 10$

$n_2 \hat{p}^* \geq 10$ $n_2(1-\hat{p}^*) \geq 10$

$$\begin{aligned}\hat{p}^* &= \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2} \\ &= \frac{109(0.60) + 91(0.65)}{109 + 91} \\ &\approx 0.62\end{aligned}$$

Under H_0

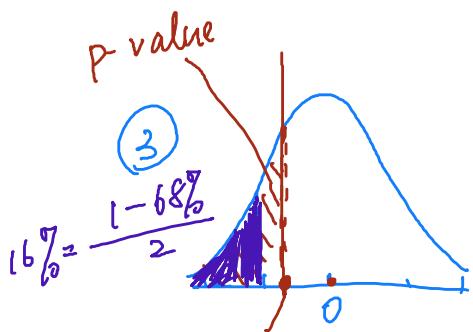
$$\hat{P}_1 \sim N(p_1, \sqrt{\frac{p_1(1-p_1)}{n_1}}) \quad X$$

$$\hat{P}_2 \sim N(p_2, \sqrt{\frac{p_2(1-p_2)}{n_2}}) \quad Y$$

$$\text{Var}(X-Y) = \text{Var}(X) + \text{Var}(Y)$$

$$\hat{P}_1 - \hat{P}_2 \stackrel{H_0}{\sim} N(0, \sqrt{\frac{\hat{p}^*(1-\hat{p}^*)}{n_1} + \frac{\hat{p}^*(1-\hat{p}^*)}{n_2}})$$

$$N(0, \sqrt{\hat{p}^*(1-\hat{p}^*)(\frac{1}{n_1} + \frac{1}{n_2})}) \quad \sqrt{(0.62)(1-0.62)(\frac{1}{109} + \frac{1}{91})} = 0.069$$



$$Z = \frac{\hat{P}_1 - \hat{P}_2 - 0}{\sqrt{\hat{p}^*(1-\hat{p}^*)(\frac{1}{n_1} + \frac{1}{n_2})}} = \frac{0.6 - 0.65}{0.069} = -\frac{0.05}{0.069} \approx -0.725$$

$$z = -0.725$$



④ p-value ≥ 0.16 using empirical rule.

⑤ Conclusion: With a large p-value = 16% $> \alpha = 5\%$, we have no strong evidence to claim trt 2 is better than trt 1.

Confidence Interval :

$$\text{for } p : \hat{p} \pm z^* \cdot SE \quad SE = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$\text{for } p_1 - p_2 : \hat{p}_1 - \hat{p}_2 \pm z^* \cdot SE \quad SE = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

$$z^* : 90\% \quad 1.645$$

$$95\% \quad 1.96$$

$$99\% \quad 2.58$$

Sample Size for a target M.O.E. of p

$$\text{M.O.E.} = z^* \cdot \sqrt{\frac{p(1-p)}{n}} \leq x$$

$$\text{if prior info } p \quad n \geq \left(\frac{z^*}{x}\right)^2 \hat{p}(1-\hat{p}) \quad \text{always round UP}$$

$$\text{if no prior info} \quad n \geq \left(\frac{z^*}{x}\right)^2 \cdot \left(\frac{1}{4}\right)$$

Chapter 7

① Hypothesis Test for μ (one sample mean)

CLT for \bar{X}

distribution.

① independent data X_1, X_2, \dots, X_n independently collected from a

② * if X_1, X_2, \dots, X_n are from a Normal distribution ✓

or

* $n \geq 30$, no extreme outliers ⚡

or

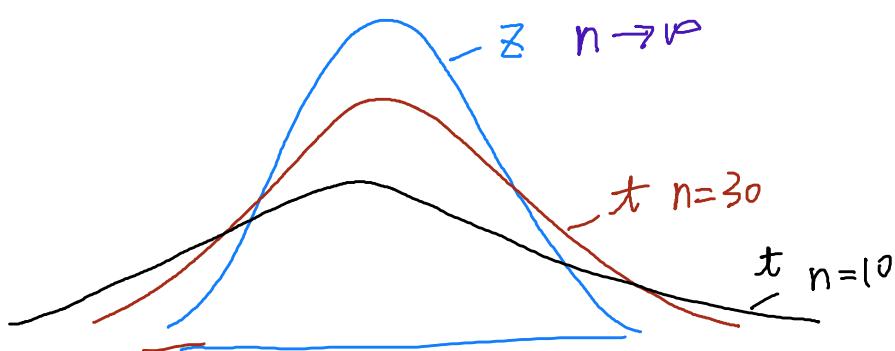
* $n < 30$, no outlier

$$\bar{X} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$$

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

$$E(\bar{X}) = (E(X_1) + E(X_2) + \dots + E(X_n)) \frac{1}{n} = \frac{1}{n} (\mu + \mu + \dots + \mu) = \mu$$

$$\text{Var}(\bar{X}) = \left(\frac{1}{n} \right)^2 (\text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n)) = \frac{1}{n^2} (\sigma^2 + \sigma^2 + \dots + \sigma^2) = \frac{\sigma^2}{n}$$



T distribution $t = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$

- { center 0
- symmetric
- bell shape
- big tail on each side.
- degrees of freedom $n-1$

$n > 30$

HT for one sample mean (one sample t test)

one population

$$\textcircled{1} \quad H_0: \mu = \underline{\mu_0}$$

$$H_a: \mu \neq \mu_0$$

$$\begin{array}{l} \mu > \mu_0 \\ \mu < \mu_0 \end{array} \quad \left. \right\}$$

\textcircled{2} Check CLT

* independent data

* normality $\left\{ \begin{array}{l} 1 \\ 2 \\ 3 \end{array} \right.$

$$\textcircled{3} \quad t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \quad df = \underline{n-1}$$

\textcircled{4} p-value or rejection region

\textcircled{5} Conclusion:

Reject H_0 with rare observation



Confidence Interval for μ

$$\boxed{\bar{x} \pm t^*(df=n-1) \cdot \frac{s}{\sqrt{n}}}$$

$$H_0: \mu = 93.3 \quad \mu (95.3, 99.3) \quad 95\% \text{ C.I.} \rightarrow \alpha = 5\% \\ \uparrow \\ \mu_0 ? \quad 99\% \text{ C.I.} \quad \alpha = 1\%$$

② HT for paired data (dependent sample)

$$H_0: \mu_{\text{diff}} = 0$$

$$H_a: \mu_{\text{diff}} \neq 0 \quad \text{one sample t-test}$$

$$\mu_{\text{diff}} < 0$$

$$\mu_{\text{diff}} > 0$$

③ HT for two population mean

$$(1) H_0: \mu_1 = \mu_2 \quad \mu_s = \mu_{ns}$$

$$H_a: \mu_1 \neq \mu_2 \quad \mu_s \neq \mu_{ns}$$

$$(2) \text{ CLT } \overline{\bar{x}_1 - \bar{x}_2} \text{ point est for } \mu_1 - \mu_2$$

✓ CLT \bar{x}_1

- independent data
- $n \geq 30$
- or
- $n < 30$
- Normal.

$$\bar{x}_1 \sim N\left(\mu_1, \frac{\sigma_1^2}{\sqrt{n_1}}\right)$$

✓ CLT \bar{x}_2

- indep
- $n \geq 30$
- $n < 30$
- Normal.

$$\bar{x}_2 \sim N\left(\mu_2, \frac{\sigma_2^2}{\sqrt{n_2}}\right)$$

✓ independent between sample 1 and sample 2

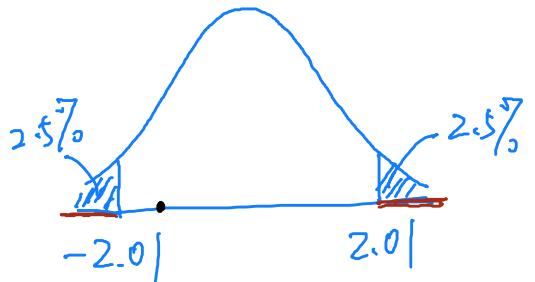
$$\bar{x}_1 - \bar{x}_2 \sim N\left(\mu_1 - \mu_2, \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right)$$

$$(3) t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{(6.78 - 7.18) - 0}{\sqrt{\frac{1.43^2}{50} + \frac{1.6^2}{100}}} = \boxed{-1.55}$$

$$df = \min(n_1 - 1, n_2 - 1) = 50 - 1 = 49$$

↓
smaller

④ Rejection Region ($\alpha=5\%$)



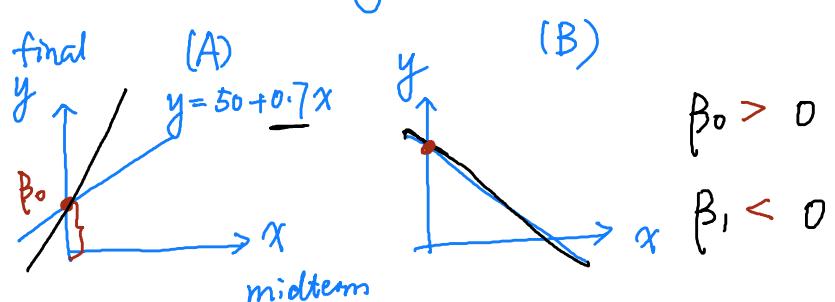
⑤ Conclusion: With test statistic -1.55 not in the rejection region

of 5% significance level, we do not have strong evidence to claim the newborns from mothers who smoke have different average birth weights than those from mothers who don't smoke.

Chapter 8 Linear Regression

Simple Linear Regression mode

$$\underline{y = \beta_0 + \beta_1 x + \epsilon}$$



x : predictor / explanatory variable / independent variable

y : outcome / response variable / dependent variable

β_0 : intercept (y value when x is zero)

β_1 : slope (change in y for a unit change in x)

Estimated Model / Fitted Line

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$\hat{y} = \underline{\underline{\beta_0 + \beta_1 x}}$$

Residuals: around zero, no special pattern

Avoid giving prediction outside available data (x) range.

Correlation Coefficient

$$R = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \cdot \frac{y_i - \bar{y}}{s_y} \right)$$

only measures linear relationship!

$$R \in [-1, 1]$$

if $|R|$ is close to 1 \rightarrow strong linear relationship

0 \rightarrow little to no linear relationship.

Final Review

Bernoulli (p)

Binomial (n, p)

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

Normal (μ, σ)

$$z = \frac{x-\mu}{\sigma} \rightarrow \text{standard normal } [0, 1]$$

p or μ ?

one p "rate" percentage / proportion " p "

$p_1 - p_2$ " " " two independent groups

one μ time / money / amount of rain --- numerical data " μ "

$\mu_1 - \mu_2$ " " " two independent groups

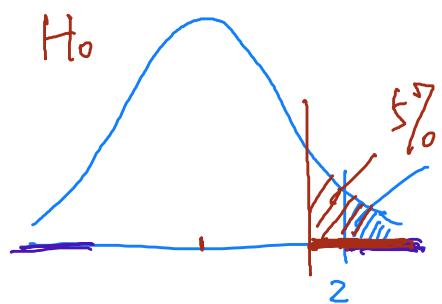
μ_{diff} difference paired data.

" σ "

Type I or Type II error
↑
(significance level)

$H_0: p = 20\%$ Youtube music subscription

$H_a: p > 20\%$ $\hat{p} = 40\%$

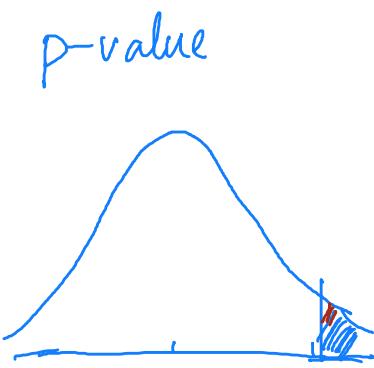


$$p\text{-value} = 0.025$$

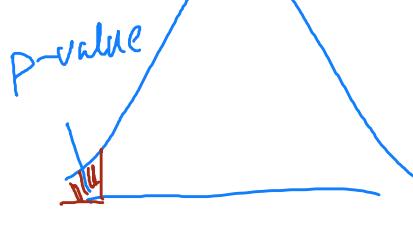
$$\alpha = 5\%$$

Reject H_0 — Type I H_0 true, reject

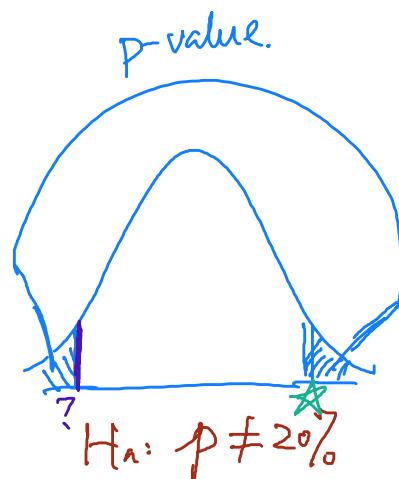
Type II H_0 false fail to reject H_0



$$H_a: p > 20\%$$



$$H_a: p < 20\%$$



$$H_a: p \neq 20\%$$

p-value: how often to observe data equally or more extreme than sample statistic.